



# Near-optimal Evasion of Randomized Convex-inducing Classifier in Adversarial Environments

PRESENTER: POORIA MADANI (@POORIA\_MDN)

FACILITATOR: TAHSEEN SHABAB

MAY 23<sup>RD</sup> , 2019



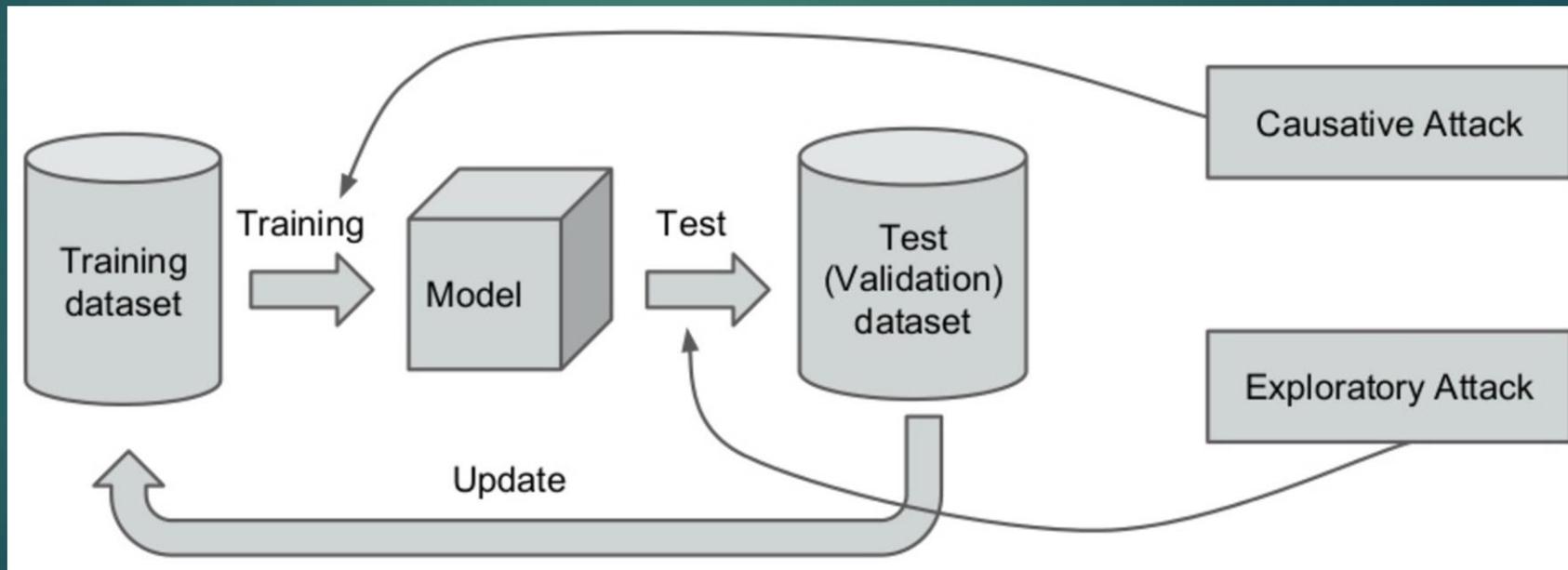
**Cybersecurity**  
is a  
game!

# Outline

- ▶ Threat Models
- ▶ Define Adversarial Classifier Reverse Engineering (ACRE)
  - ▶ What is an optimal evasion?
- ▶ Optimal Evasion of Convex Inducing Classifiers
  - ▶ Complexity Analysis
  - ▶ Moving Target Defense (MTD)
- ▶ Active Line Search
- ▶ Simulations and Results
- ▶ Discussions

# Motivation: Threat Models

- ▶ Adversaries are **smart** and **motivated** individuals – they break all the **assumptions** you have made about the data.



Source: [Xiao Huang, 2015]

# Threat Models – cont'd

	<b>Integrity</b>	<b>Availability</b>	<b>Confidentiality</b>
<b>Run-time</b>	Evasion	-	Model Extraction (e.g., hill climbing attack)
<b>Training-time</b>	Poisoning (i.e., implementing backdoors)	Poisoning (i.e., maximizing classification error)	-

# Threat Models – cont'd

## ▶ Optimal Evasion

*evade detection with minimum number of change made to an attack instance/vector.*



Source: [Biggio et. al., ACML 2011]

# Threat Models – cont'd

## ▶ Problem Setup

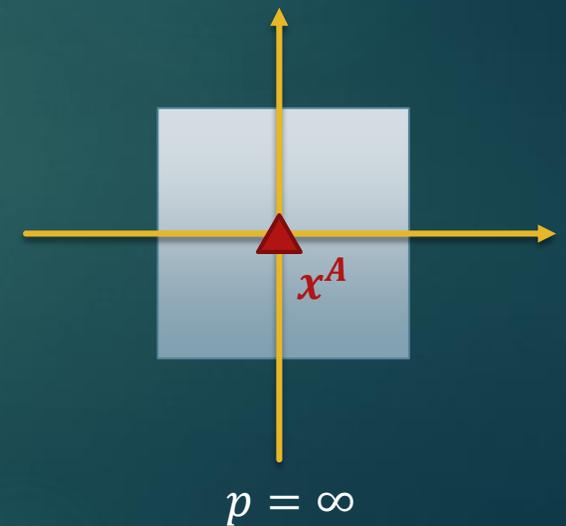
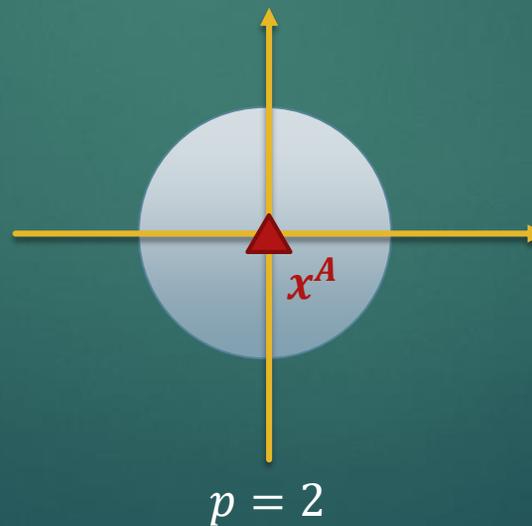
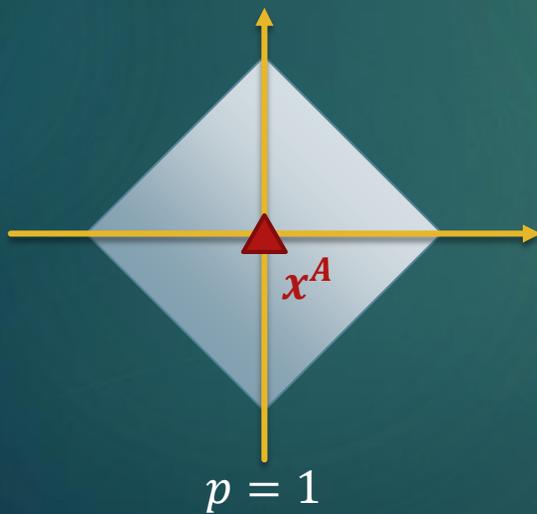
- ▶ Adversary has a detectable malicious instance  $x^A$
- ▶ Adversary is aware of one instance  $x^-$  that belongs to benign class
- ▶ Adversary can construct an arbitrary  $x$  and query the detection system to find out its label (e.g., +/-)



# Minimal Adversarial Cost (MAC)

- ▶ Adversarial Cost Function

$$A_p^c(x - x^A) = \left( \sum_{d=1}^D c_d |x_d - x_d^A|^p \right)^{1/p}$$



# Minimal Adversarial Cost (MAC) – cont'd

9

- ▶ Adversarial Cost Function

$$A_p^c(x - x^A) = \left( \sum_{d=1}^D c_d |x_d - x_d^A|^p \right)^{1/p}$$

Initial Malicious Instance

- ▶ Then, for a given classifier  $f$

$$MAC(f, A) \triangleq \inf_{x \in X_f^-} [A(x - x^A)]$$

- ▶ Instance of minimal adversarial cost

$$\epsilon - IMAC(f, A) \triangleq \{x \in X_f^- \mid A(x - x^A) \leq (1 + \epsilon) \cdot MAC(f, A)\}$$

Benign Class

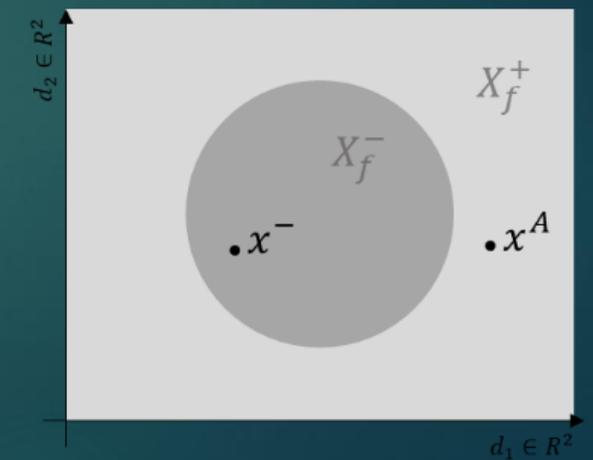
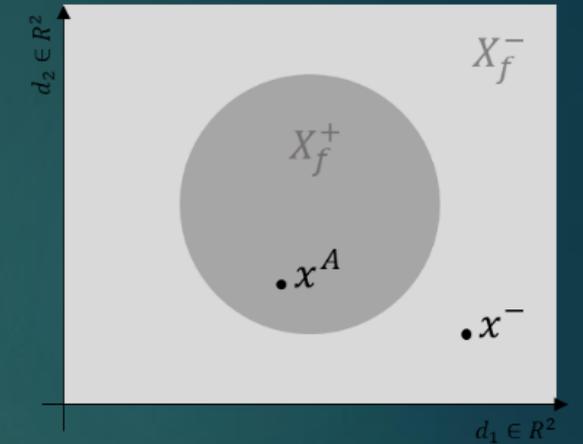
- ▶ Near Optimal Evasion (ACRE searchable)

A family classifiers  $F$  is  $\epsilon - IMAC$  searchable under a family cost function if for all  $f \in F$  and some  $A$ , there is an algorithm that finds  $x \in \epsilon - IMAC$  using polynomial-many membership queries. [Nelson et al. 2010]

# ACRE Searchability of Convex Bodies

10

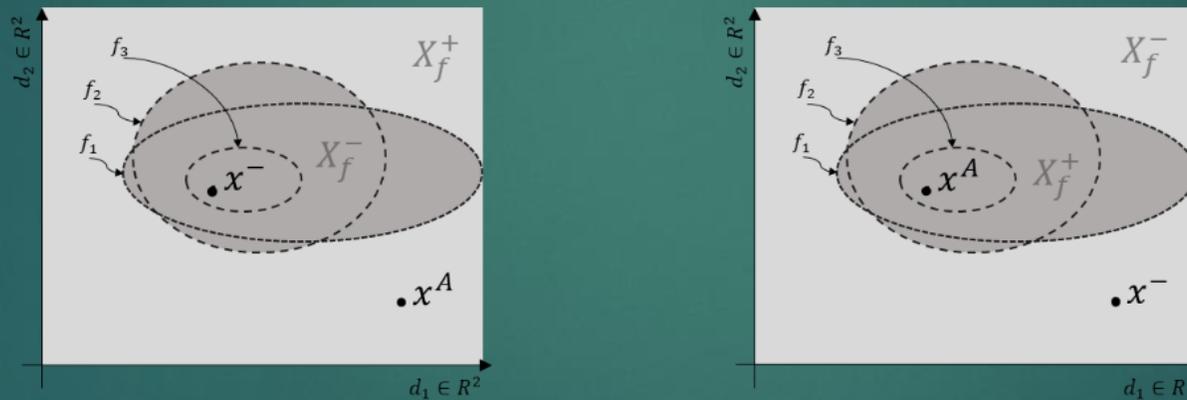
- ▶ Convex-family classifiers forms an important collection of algorithms used in many adversarial environments.
  - ▶ *linear classifiers, hyper-sphere boundary anomaly detectors, one-class anomaly detectors, and ...*
- ▶ For  $\ell_1$  (i.e.,  $p=1$ ) adversarial cost function, the family of convex inducing classifiers are ACRE searchable.
- ▶ For  $\ell_{>1}$  the problem is NP-Complete.



# Moving Target Defense

11

- ▶ It is **hypothesized**, that use of multiple classifiers at runtime and randomly selecting among them at runtime could make ACRE searchability harder.



- ▶ Construct multiple models, select one at random to handle a given query.

# Reformulation of ACRE Searchability

12

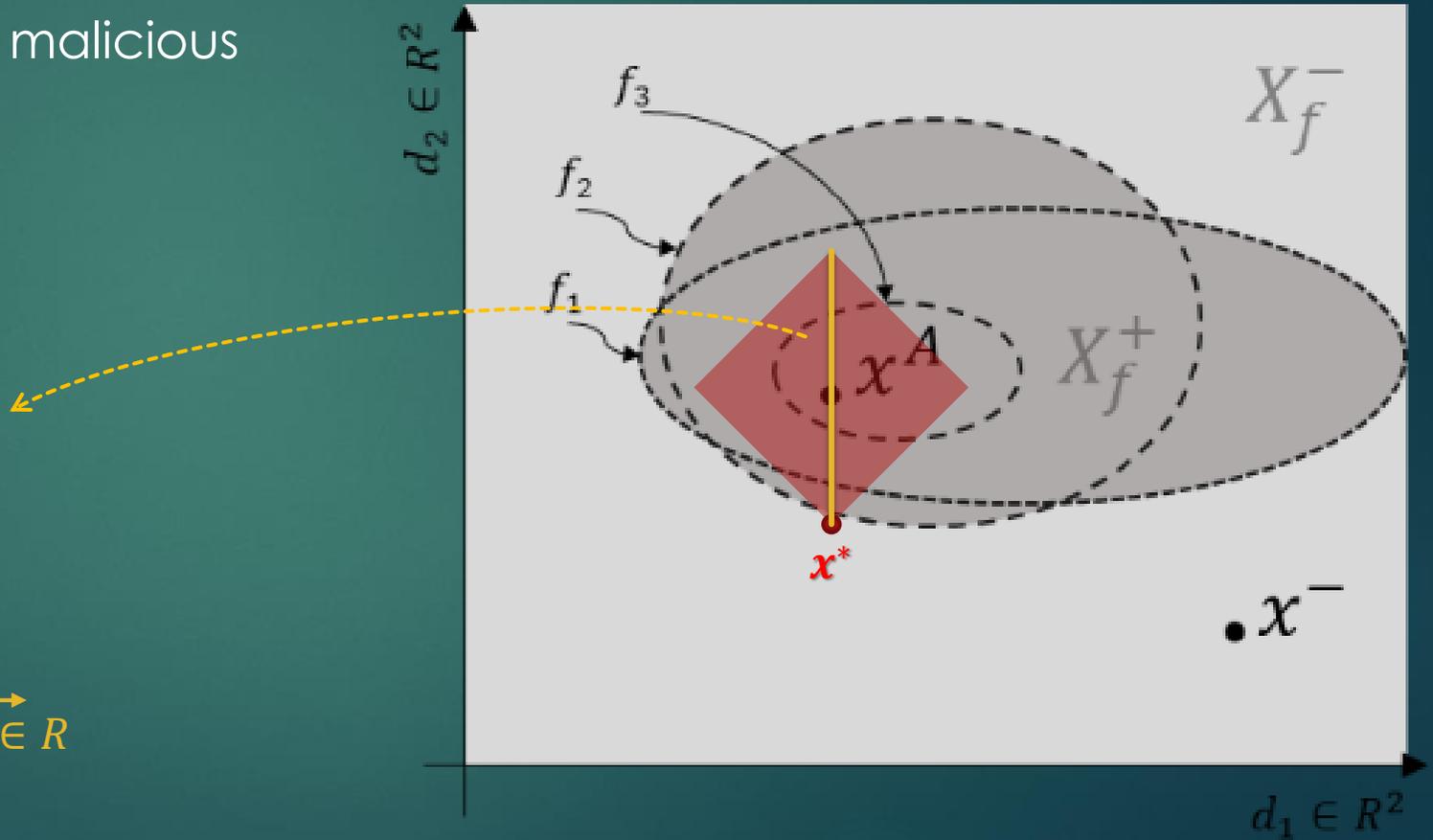
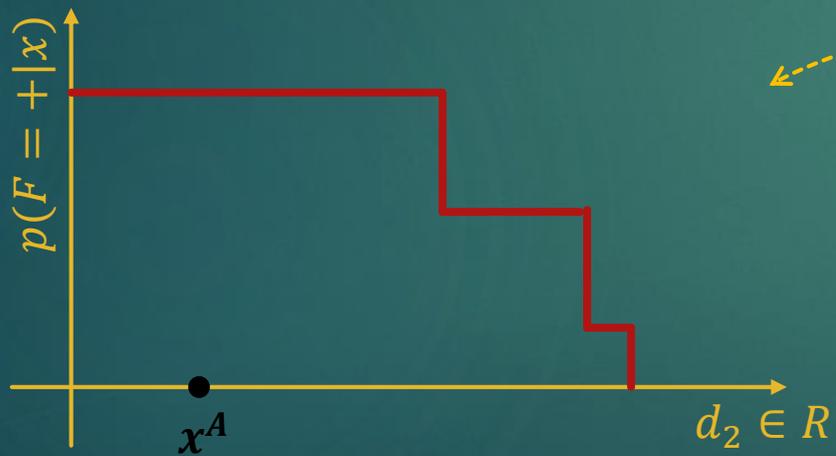
- ▶ Given the randomized strategy of the defender:  
Find an instance of  $\epsilon$ -IMAC that evade detection with probability better than  $P$ .
- ▶ The goal of adversary is construct



$$\text{w.r.t. } \inf_{x^* \in X_f^-} A_1^c(x^* - x^A)$$
$$p(f = +|x^*) \leq P$$

# Active Line Search

- ▶ For  $\ell_1$  cost function and a malicious (i.e.,  $X_f^+$ ) convex family



# Active Line Search – cont'd

14

- ▶  $p(F = +|x, d)$  is representing the probability of success of a Bernoulli random variable at point  $x$ .
- ▶ Use Gaussian Process Regression (GPR) to learn  $p(F = +|x, d)$  for different  $x$  along each feature dimension.
  - ▶ Learn the probability of success ( $\theta$ ) through sampling and submitting **membership queries** to the defender's black box.

$$\begin{bmatrix} \theta \\ \theta^* \end{bmatrix} \sim GPR(\mu, \Sigma) \quad \Sigma = \begin{pmatrix} \kappa + \sigma(\theta)I & k_* \\ k_*^T & k_{**} \end{pmatrix} \quad \kappa(q_i, q_j) = \exp\left(-\frac{1}{2l^2} (q_i - q_j)^2\right)$$

- ▶ Interpolate  $\theta^*$  for other  $x$  that have not been observed.

$$\theta^* = k_*^T (\kappa + \sigma I)^{-1} \theta \quad \text{Var}[\theta^*] = k_* - k_*^T (\kappa + \sigma I)^{-1} k_*$$

# Active Line Search – cont'd

Algorithm 1: Algorithm to find initializer.

```
Data:  $x^A, x^-, C, \epsilon, T_\sigma, P$   
Result:  $x^*$   $p\epsilon$ -IMAC instance  
1  $\vec{\theta} \leftarrow [1, 0]$   
2  $\vec{X} \leftarrow [x^A, x^-]$   
3  $x^* \leftarrow x^-$   
4 do  
5   construct GPR (MEAN( $\vec{\theta}$ ), KERNEL( $\vec{X}$ ))  
6   foreach  $d_i \in D$  do  
7     foreach  $\{x | x^A \pm C_d \cdot \delta_d \text{ that are } \epsilon \text{ spaced}\}$  do  
8        $\theta_i, var_i \leftarrow \hat{f}(x^i)$ ; // GPR prediction at  $x^i$   
9       if  $var_i$  is maximum then  
10        |  $max_i \leftarrow x^i$   
11        end  
12        if  $\theta_i < P$  and  $var_i \approx 0$  and  $|x_d^A - x_d^i| \leq C_d$   
13          then  
14            |  $C_d \leftarrow |x_d^A - x_d^i|$   
15            | if  $\|x^A - x^i\|_1 \leq \|x^A - x^*\|_1$  then  
16              |  $x^* \leftarrow x^i$   
17            | end  
18          end  
19        end  
20      query the oracle and find  $\theta_i$  and  $\sigma(\theta_i)$  for  $max_i$   
21       $\vec{X} \leftarrow \vec{X} \cup max_i$   
22       $\vec{\theta} \leftarrow \vec{\theta} \cup \theta_i$   
23      update KERNEL and noise  $\sigma(\theta_i)$  based on Equation 4  
24 while maximum variance is above  $T_\sigma$ ;  
25 return  $x^*$ 
```

▶ Begin by two known points to the adversary:  $x^A$  and  $x^-$  and their corresponding  $\theta = 1$  and  $\theta = 0$  and train a GRP.

▶ Then, get the trained GRP to interpolate all ( $\epsilon$  separated) data points along each feature vector, and select the one with **highest variance** ( $x_i$ ).

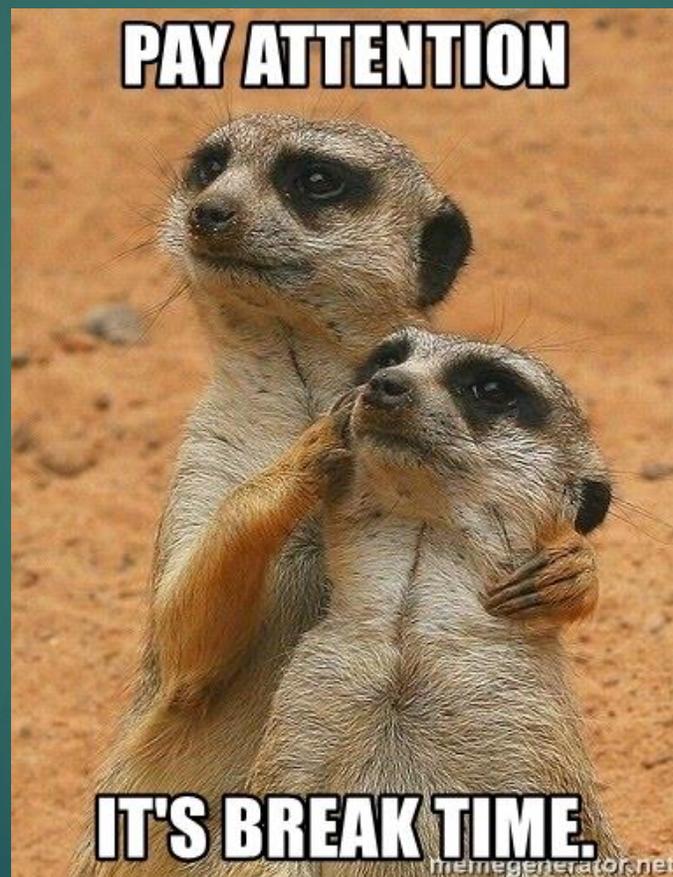
▶ Find an estimate for  $\theta_i$  of the selected point by multiple membership queries.

▶ Update the covariance matrix with newly acquired information and **uncertainty**.

▶ Retrain the GRP

$$Var(\theta_i) = \sqrt{\frac{\theta_i(1-\theta_i)}{N}}$$

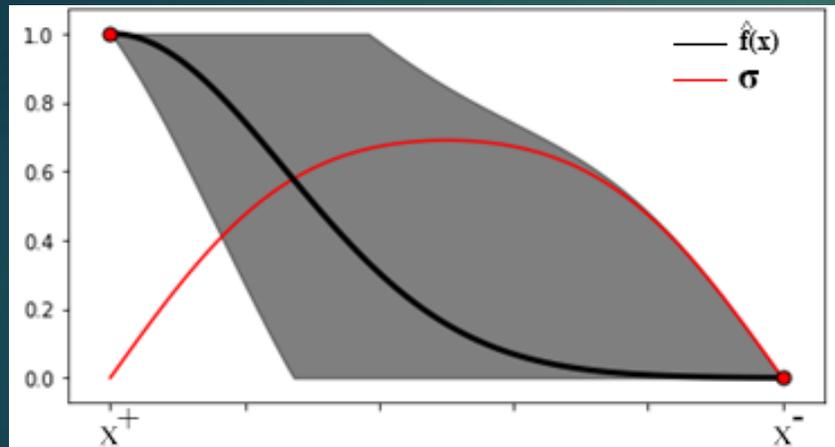
# BREAK



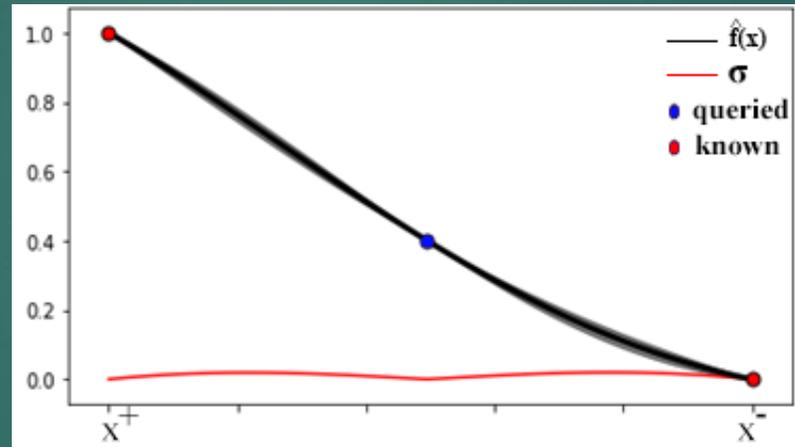
# Toy Example (1)

- ▶ Example of 1-D randomized decision points

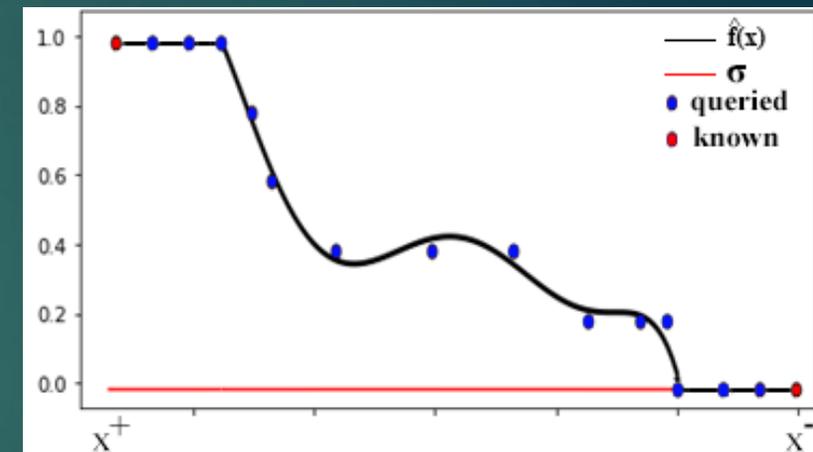
Initialization



Iteration 1

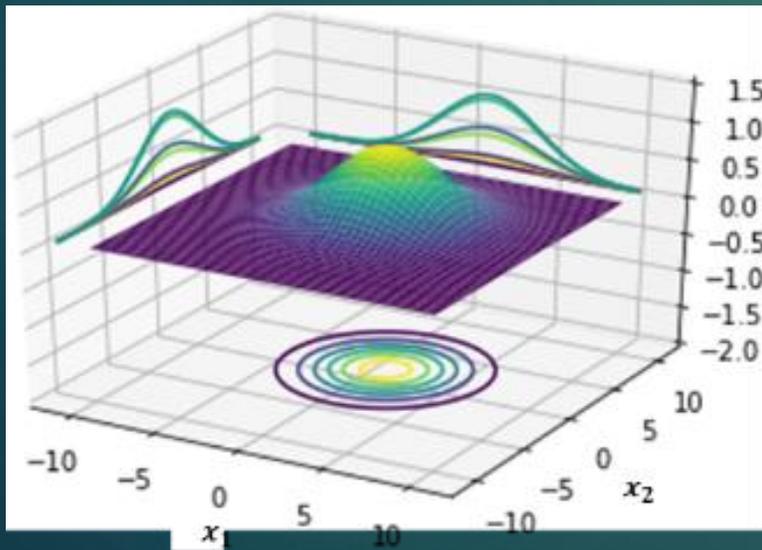
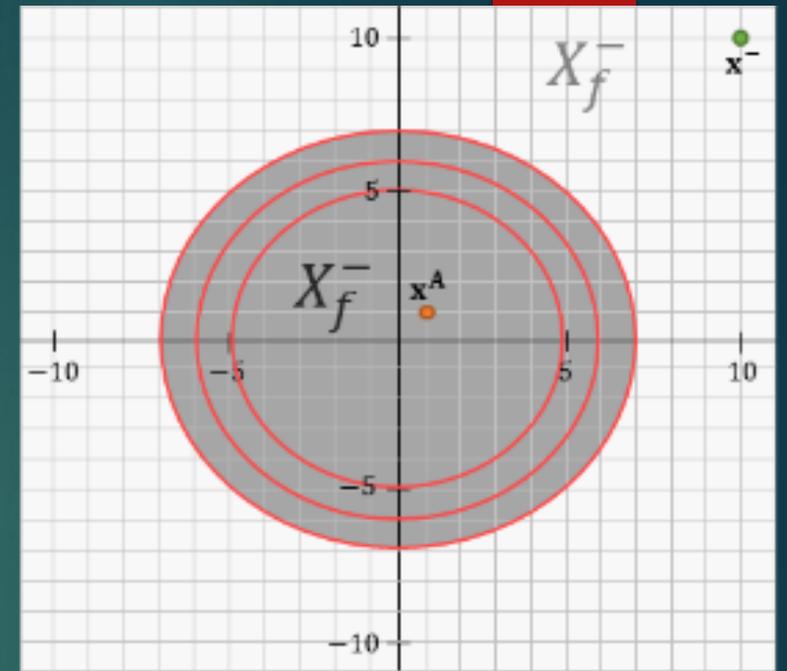


Last Iteration

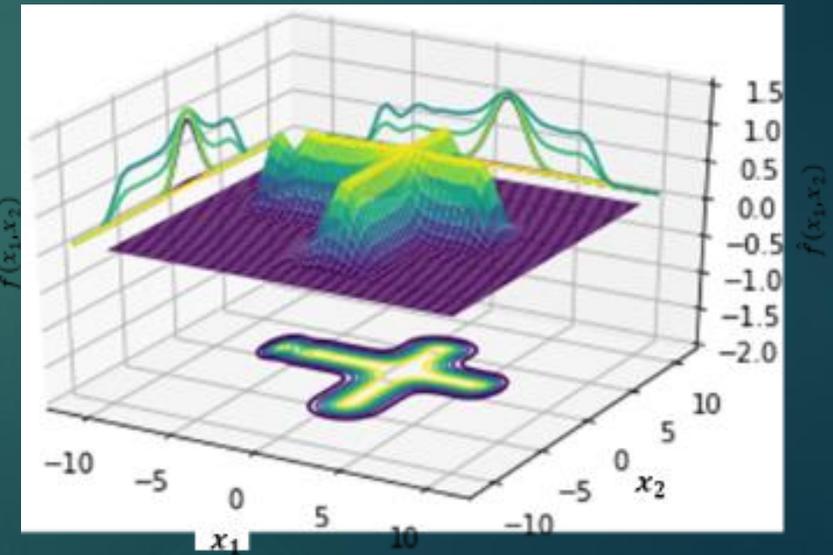
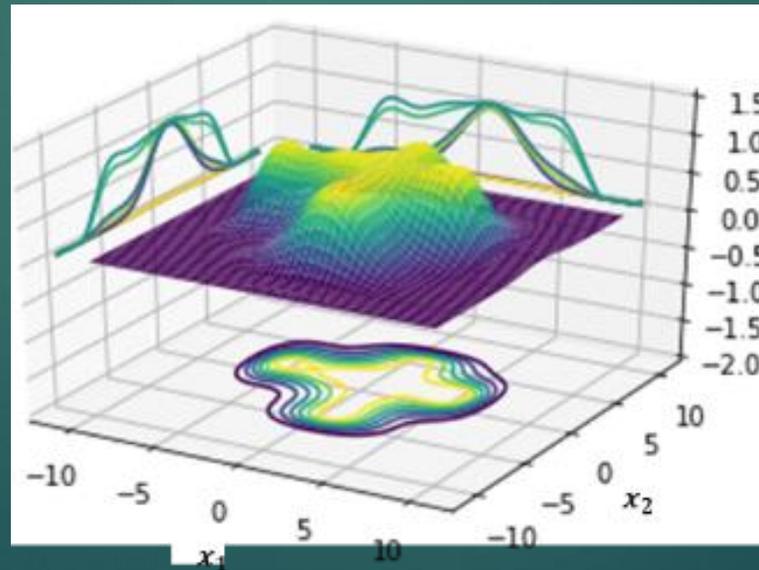


# Toy Example (2)

- ▶ Example of 2-D randomized decision hyper-spheres



Initialization



Final Iteration

# Toy Examples' Analysis

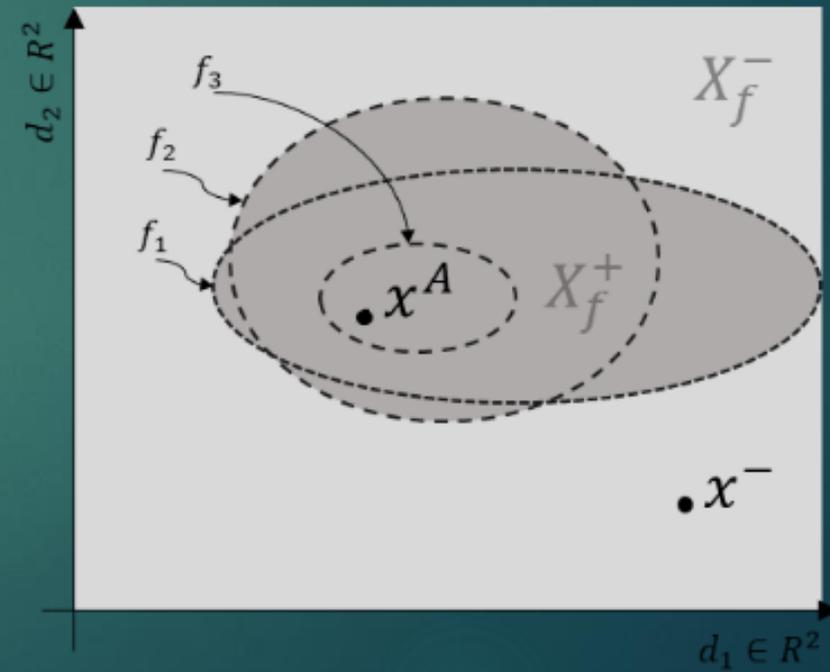
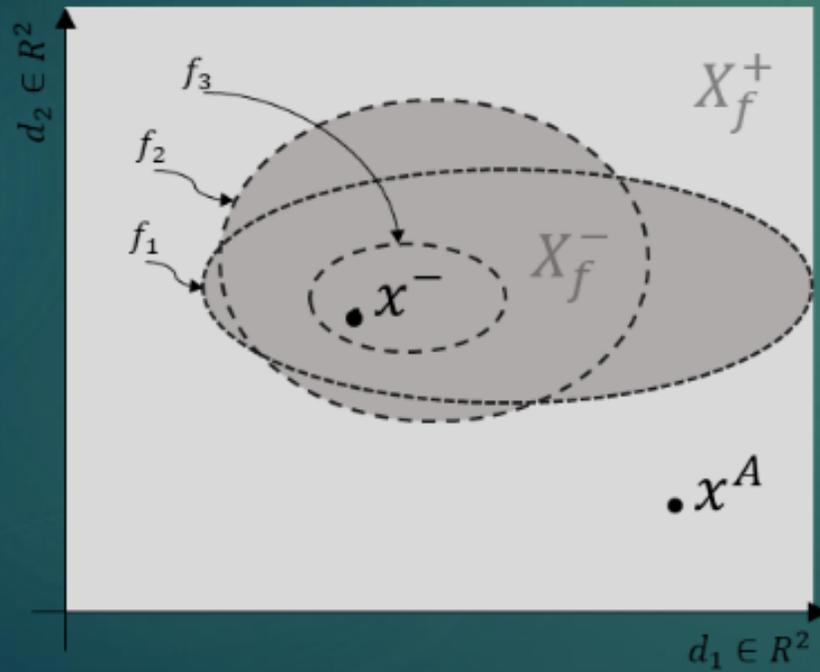
- ▶ Randomization certainly have increased the number of required query.
- ▶ Ignoring randomization can result in non-optimal result.
- ▶ Our proposed active line search algorithm converges with time complexity of

$$O\left(d \cdot \frac{|x^A - x^-|}{\epsilon}\right).$$

Method	To Example	P	#Queries	Avg Error relative to MAC
Nelson et al.	1-D	1.0	<b>34</b>	-0.6183020
		0.5	N/A	N/A
	2-D	1.0	<b>57</b>	-1.5764189
		0.5	N/A	N/A
Active Multiline Search	1-D	1.0	122	<b>0.0022201</b>
		0.5	122	<b>0.0175023</b>
	2-D	1.0	254	<b>0.5510204</b>
		0.5	254	<b>0.5306122</b>

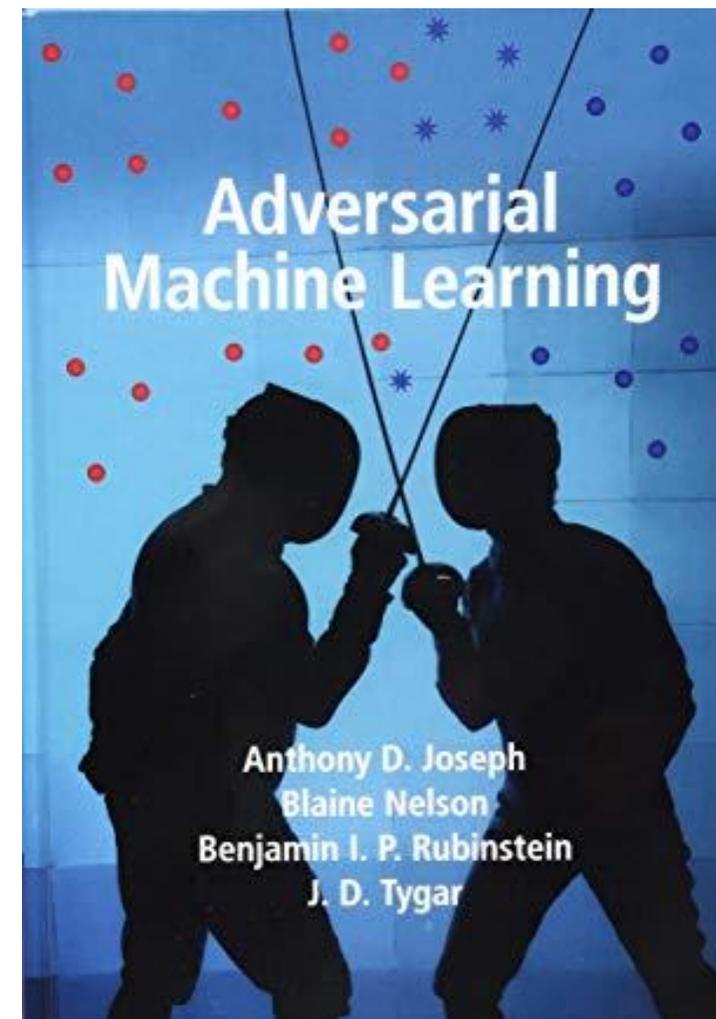
# What about when $X_f^-$ is convex?

- ▶ We believe such case is not ACRE searchable!  
Can you help us prove it?



# Closing Remarks

- ▶ Randomization, certainly, would create more trouble for an adversary to successfully carry-out an ACRE search.
- ▶ The randomization, when  $X_f^+$  is convex, can still be evaded using polynomial-many queries.
- ▶ How would randomization affect ACRE search when  $X_f^-$ ?
- ▶ How would feature-space transformation (with or without) randomization affect ACRE search?



# Discussion Points

- ▶ Does randomization have an opportunity cost?
- ▶ I can't help but think about the early misconceptions of encryption where randomized hashing is/was often termed as more secure by novice practitioners which is not as per best practice. Is this true for randomization in this project?
- ▶ Would you happen to know if there are controls in place today to detect such abnormal sequence of inputs from potentials hackers?
- ▶ Would it be possible to elaborate on why one is a convex problem and not the other?

# Discussion Points – cont'd

23

- ▶ Have you experimented with the constant which bounds the binomial sampling?
  - ▶ How do we pick the optimal constant?
- ▶ Do you think the sampling performed by Hackers can trigger existing anomaly detection tools used in production environments, thereby preventing the attack altogether.
  - ▶ Ex. If we were to use the above approach on GMAIL, do you think google could possibly model the sampling (ex. Sending the same email with small deviations in an iterative manner) as an anomaly? STOP THE PROBLEM AT THE SOURCE
- ▶ In our data-driven world, do you think it makes sense for SOC Analysts to have a visualization of decision boundaries handy for them to audit?
  - ▶ Rule-based detection is audited by experts for IPS/IDS.
  - ▶ What do you think is a best practice for us going forward?

# Discussion Points – cont'd

24

- ▶ Would you have any intuition on explainability of models leveraging the MAB model you depicted?