# On the Hardness of Robust Classification

Pascale Gourdeau, Varun Kanade, Marta Kwiatkowska and James Worrell

Department of Computer Science, University of Oxford

## Question

*What distributional assumptions are needed and how much power can we give an adversary to ensure efficient robust learning?*

## Take Away

- Inadequacies of widely-used definitions of robustness surface under a learning theory perspective.
- It may be possible to only solve robust learning problems with strong *distributional assumptions*.
- Simple proof for computational hardness of robust learning.

## Problem Setting

Our paper:

- Binary classification
- Binary feature vectors (input space: $\mathcal{X} = \{0,1\}^n$)
- An adversary can modify input bits after training (*evasion attacks*)

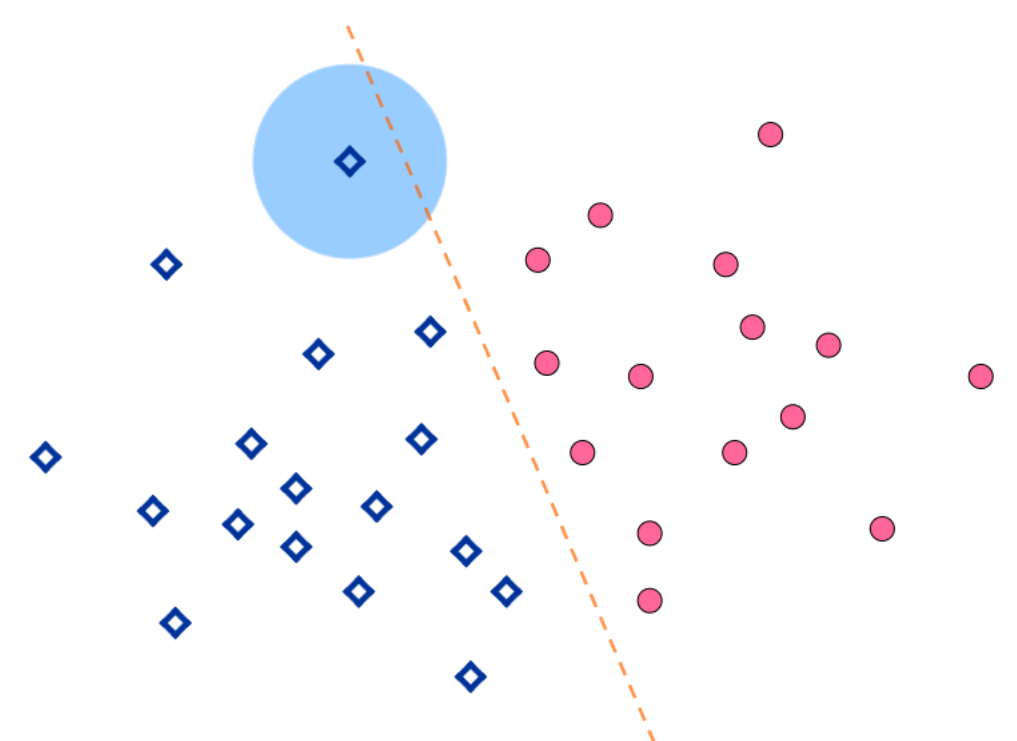For example, we wish to be able to differentiate between 0's and 1's:



The image of a 0 should not be classified as a 1 if it is slightly perturbed by an adversary:



**Efficient Robust Learning:**

In general, we want to prove or disprove the existence of an algorithm with *polynomial sample complexity* (in the learning parameters and input dimension $n$) that will output a hypothesis such that the probability of drawing a new point that can be perturbed by an adversary and resulting in a misclassification to be small:



*But what counts as a misclassification?*

## Robust Risk Definitions

**Constant-in-the-ball:**

$$\mathsf{R}_\rho^C(h,c) = \mathbb{P}_{x\sim D}\left(\exists z \in B_\rho(x) : h(z) \neq c(x)\right)$$
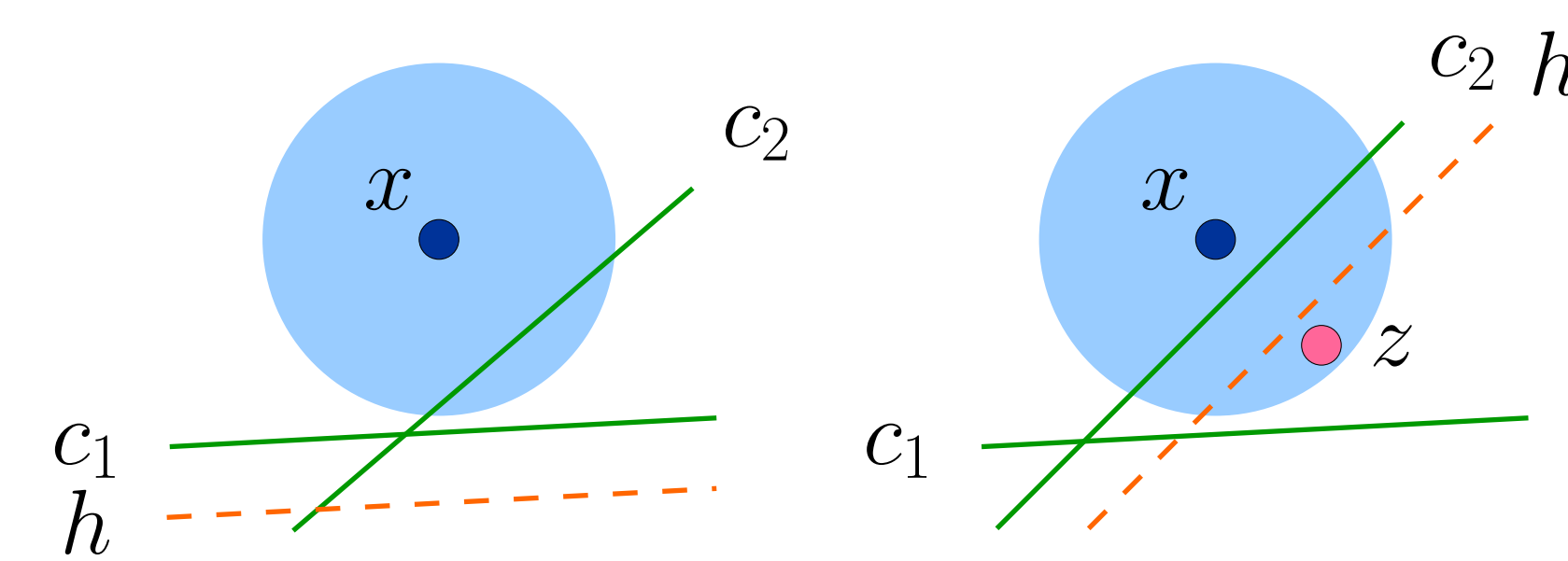


Figure: The robust loss is 0 on the LHS and 1 on the RHS.

**Exact-in-the-ball:**

$$\mathsf{R}_\rho^E(h,c) = \mathbb{P}_{x\sim D}\left(\exists z \in B_\rho(x) : h(z) \neq c(z)\right)$$
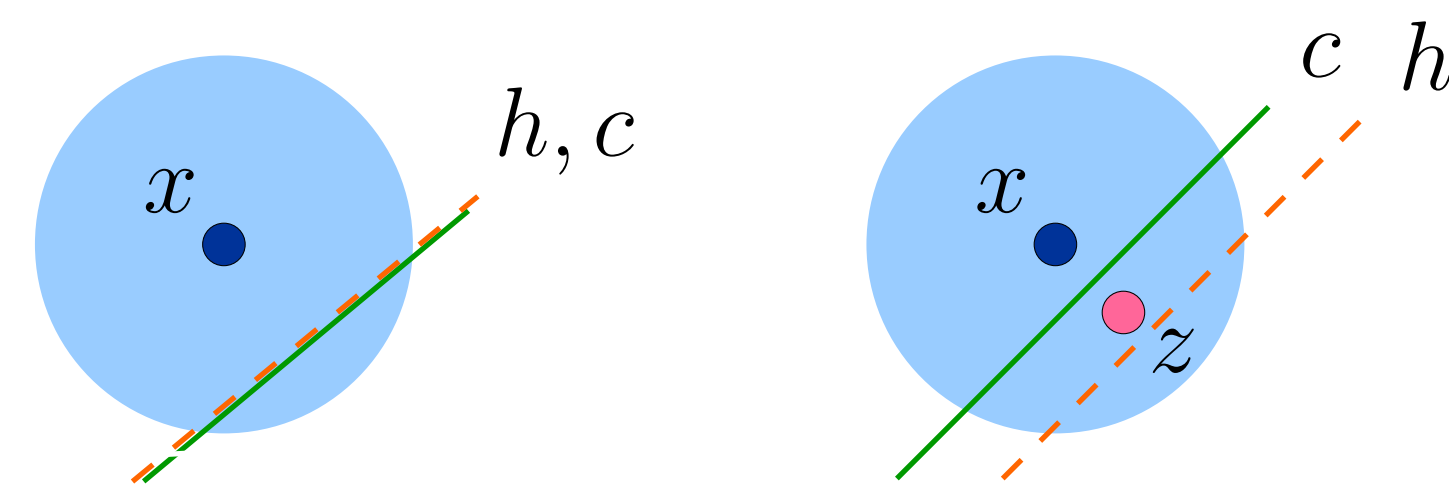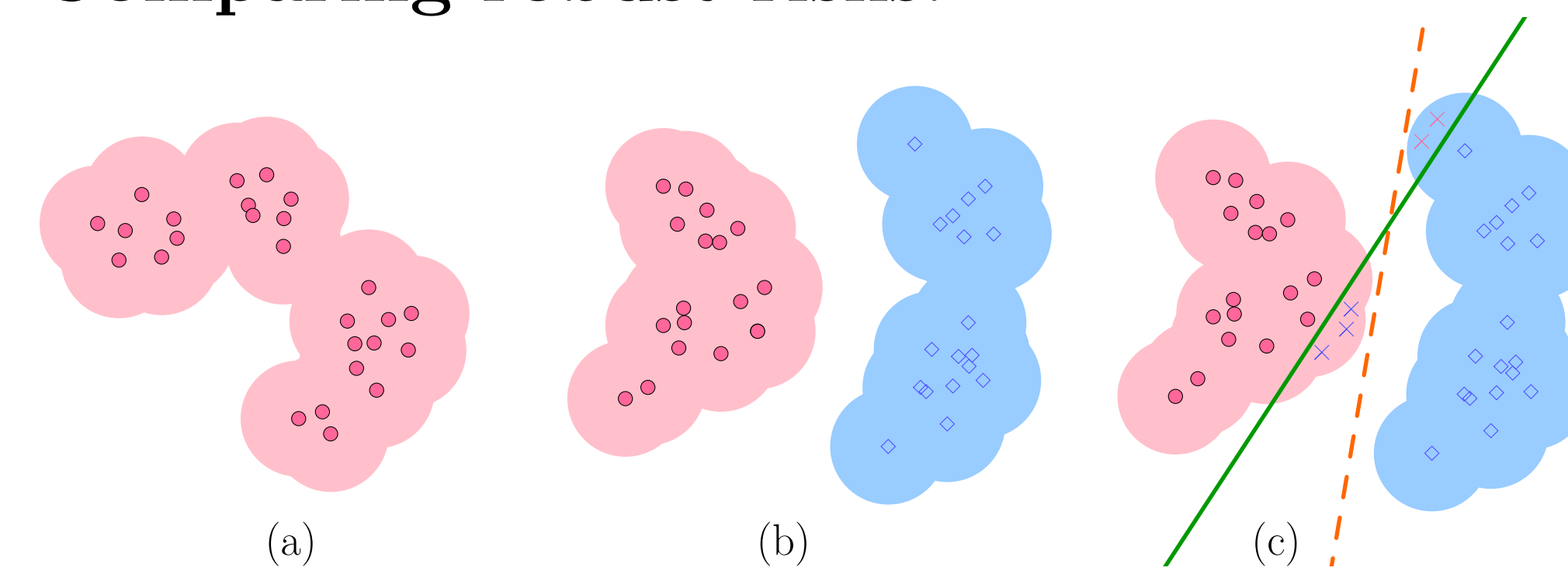


Figure: The robust loss is 0 on the LHS and 1 on the RHS.

**Comparing robust risks:**



(a) $\mathsf{R}_\rho^C(h,c) = 0$ only achievable if $c$ is constant.
(b) There exist $h$ such that $\mathsf{R}_\rho^C(h,c) = 0$.
(c) $\mathsf{R}_\rho^C$ and $\mathsf{R}_\rho^E$ differ. The solid concept is the target, while the dashed one is the hypothesis. Shaded regions represent the dots' $\rho$-expansion. The crosses are perturbed inputs causing $\mathsf{R}_\rho^E > 0$, while $\mathsf{R}_\rho^C = 0$.
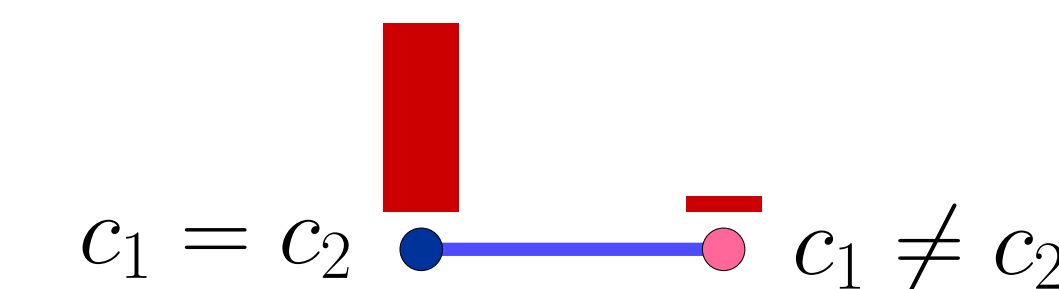
To us, the adversary's power: creating perturbations that cause the hypothesis and target functions to disagree, so we use the *exact-in-the-ball* definition.

## Distribution-Free Robust Learning

**Theorem:** *Any concept class $\mathcal{C}$ is efficiently distribution-free robustly learnable if and only if it is trivial.*

A class of functions is *trivial* if $\mathcal{C}_n$ has at most two functions, and that they differ on every point.

Distributional assumptions are *essential*:



$c_1 = c_2$ — $c_1 \neq c_2$

## Monotone Conjunctions

**Question:** How much power can we give an adversary and still ensure efficient robust learnability?

**Monotone conjunctions:**

thesis $\wedge$ sleep deprivation $\wedge$ caffeine

**Theorem:** *The threshold to robustly learn monotone conjunctions under log-Lipschitz distributions is $\rho(n) = O(\log n)$.*

$\rho = O(\log n)$: PAC algorithm is a robust learner.
$\rho = \omega(\log n)$: no sample-efficient learning algorithm exists.

**$\alpha$-Log-Lipschitz Distributions:**

$$x_1 = (0,\ldots,1,1,1,\ldots,0)$$
$$x_2 = (0,\ldots,1,0,1,\ldots,0) \implies \frac{p(x_1)}{p(x_2)} \leq \alpha .$$

For e.g.: uniform distribution, product distribution where the mean of each variable is bounded, etc.

**Intuition:** input points that are close to each other cannot have vastly different probability masses.

## Computational Hardness

- An information-theoretically easy problem can be computationally hard.
- We give a simple proof of the computational hardness of robust learning result of [1].
- We reduce a computationally hard PAC learning problem to a robust learning problem.
- We use the trick from [1] of encoding a point's label in the input for the robust learning problem.

**Reduction.** Take a PAC learning problem for concept and distribution classes $\mathcal{C}$ and $\mathcal{D}$ defined on $\mathcal{X} = \{0,1\}^n$. Define $\varphi_k$ as follows:

$$\varphi_k(x) := \underbrace{x_1 \ldots x_1 x_2 \ldots x_{d-1} x_d \ldots x_d}_{2k+1 \text{ copies of each } x_i} c(x) ,$$

❶ Blow up input space to $\mathcal{X}' = \{0,1\}^{(2k+1)n+1}$.
❷ New concept class:

$$\mathcal{C}' = \{c \circ \mathrm{maj}_{2k+1} \mid c \in \mathcal{C}\} ,$$

where $\mathrm{maj}_l$ returns the majority vote on each subsequent block of $l$ bits, and ignores the last bit.
❸ Distribution family $\mathcal{D}'$: for each $c \in \mathcal{C}$, $D \in \mathcal{D}$, we have a new $D'$ as follows for $z \in \mathcal{X}'$:

$$D(z) = \begin{cases} D(x) & z = \varphi_k(x), \\ 0 & \text{otherwise.} \end{cases}$$

**Reasoning.**

- Any algorithm for learning $\mathcal{C}$ w.r.t. $\mathcal{D}$ yields an algorithm for learning the pairs $\{(c', D')\}$.
- A *robust* learner cannot only rely on the last bit of $\varphi_k(x)$ (it could be flipped by an adversary).
- A *robust* learner can be used to PAC-learn $\mathcal{C}_n$.

## References

[1] S. Bubeck, E. Price, and I. Razenshteyn. Adversarial examples from computational constraints. *arXiv preprint arXiv:1805.10204*, 2018.

[2] D. Diochnos, S. Mahloujifar, and M. Mahmoody. Adversarial risk and robustness: General definitions and implications for the uniform distribution. In *Advances in Neural Information Processing Systems*, pages 10359–10368, 2018.