

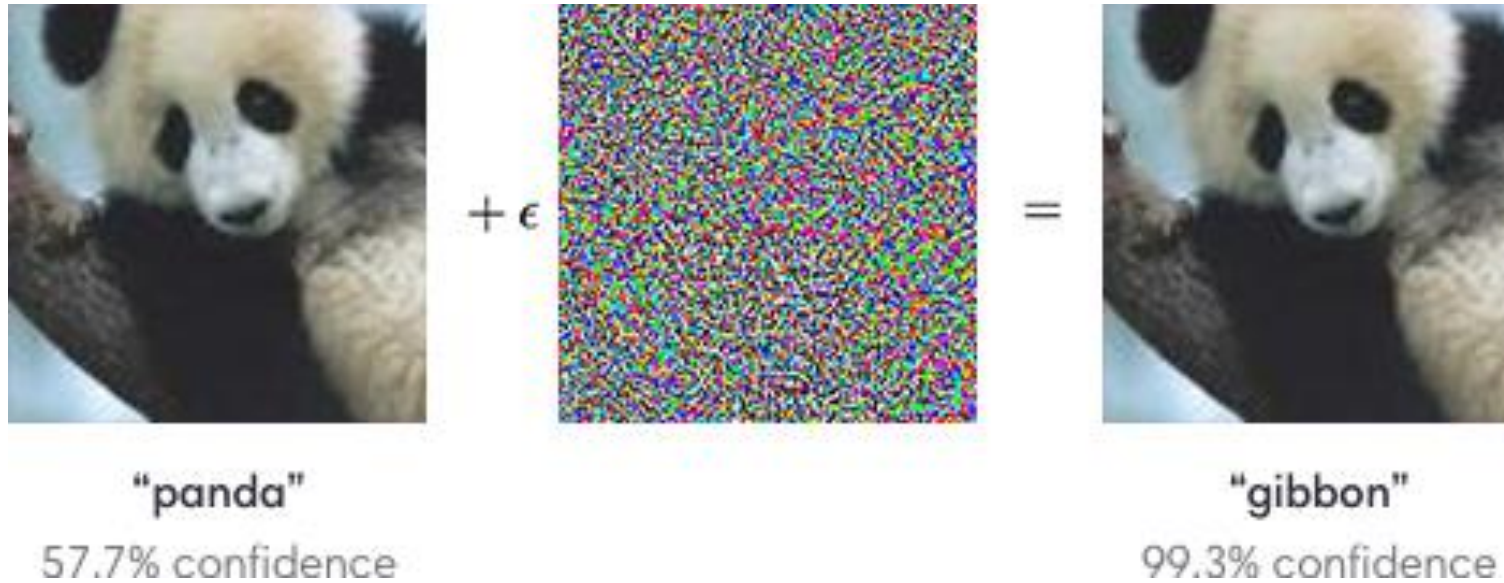
# Adversarial Robustness From Well-Separated Data

Cyrus Rashtchian  
UC San Diego

Joint with Yao-Yuan Yang, Yizhen Wang, Kamalika Chaudhuri **AISTATS 2020**

+ Hongyang Zhang, Ruslan Salakhutdinov **NeurIPS 2020**

# Adversarial Examples



- Lowd and Meek 2005
- Szegedy et al 2013
- Papernot et al 2014,15,16
- ...

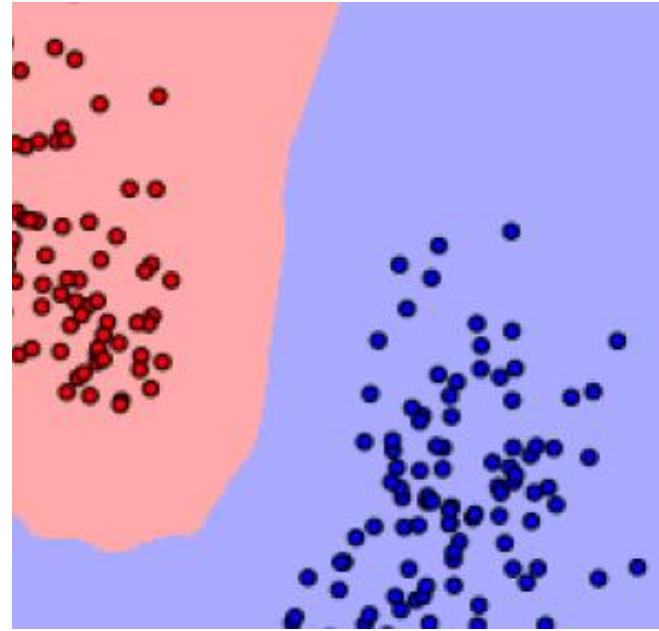


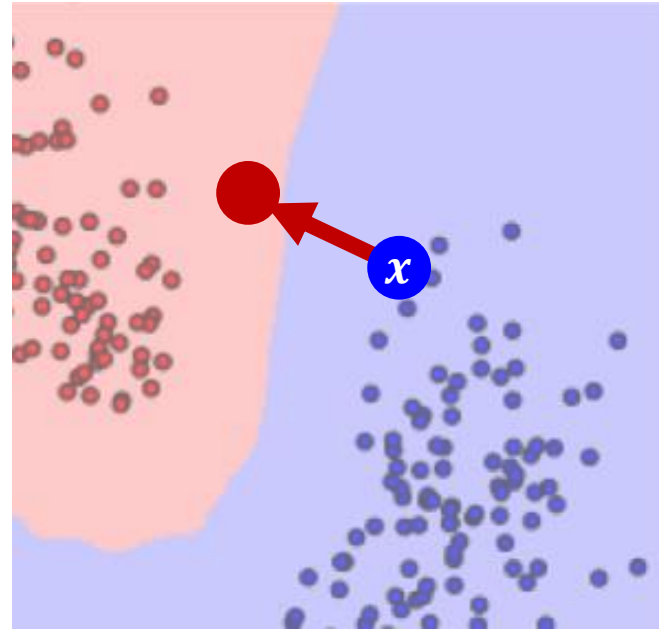
labradoodle or fried chicken



chihuahua or muffin



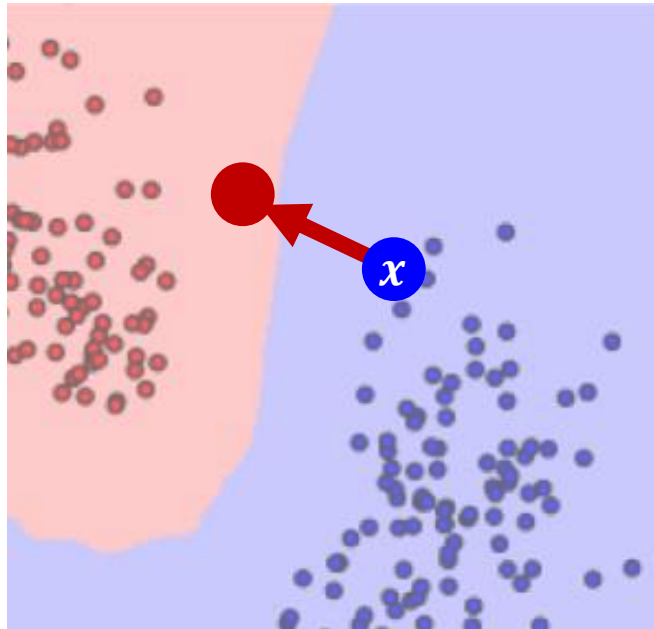




## Definition (Adversarial example)

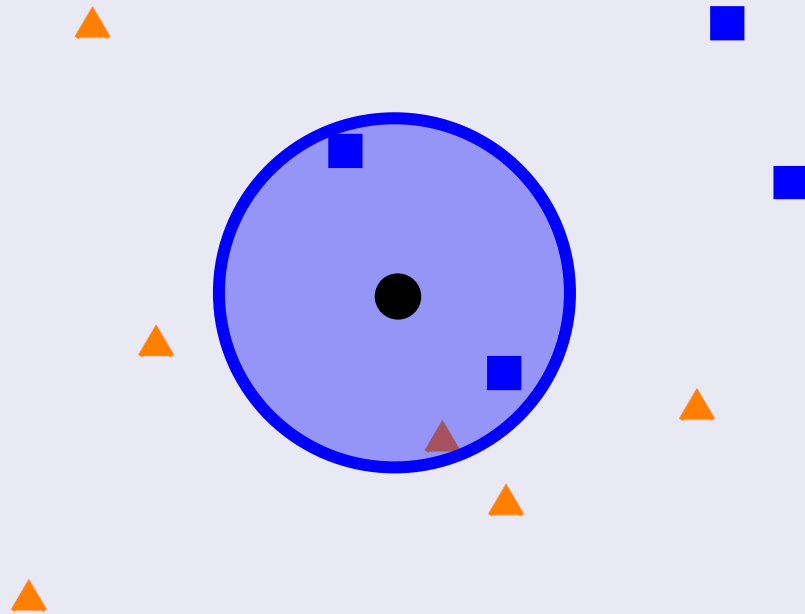
$\mathbf{x}_{adv}$  is an adversarial example of the target example  $\mathbf{x}$  if and only if

$$\|\mathbf{x} - \mathbf{x}_{adv}\|_p \leq r \text{ and } f(\mathbf{x}) \neq f(\mathbf{x}_{adv})$$



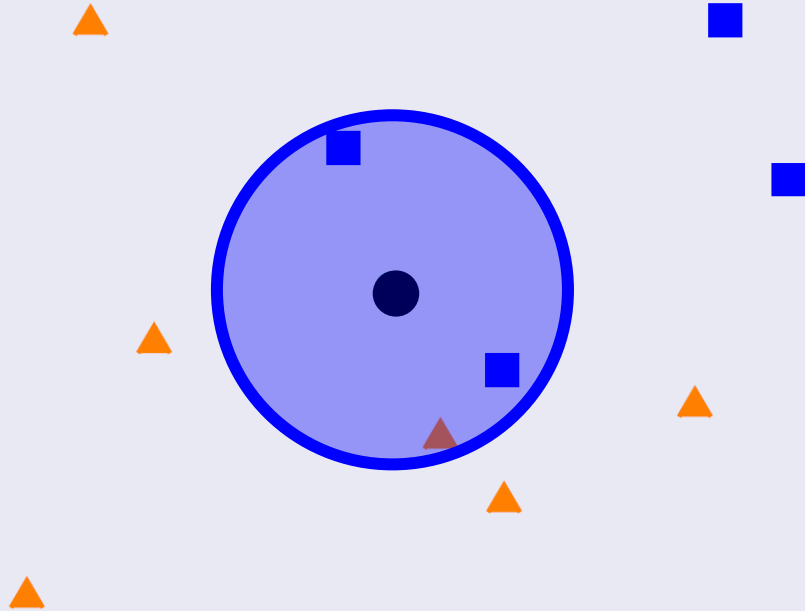
## $k$ nearest neighbor ( $k$ -NN)

take  $k$  closest training examples and output the majority label



## $k$ nearest neighbor ( $k$ -NN)

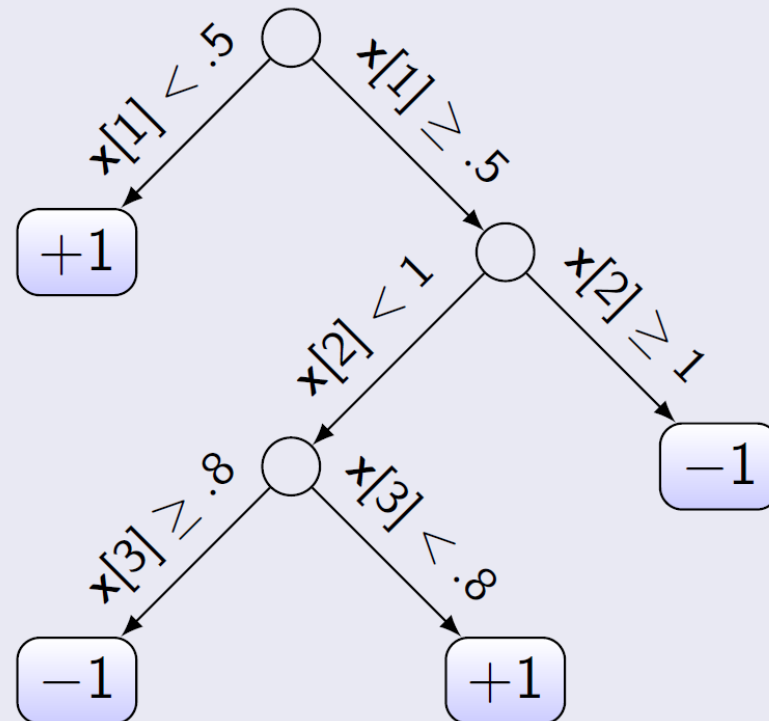
take  $k$  closest training examples and output the majority label



## Decision tree and tree ensembles

recursively split the data

- common models: decision tree, random forest, gradient boosting trees, etc.





## Optimal attack

$$\min_{\mathbf{x}_{adv}} \|\mathbf{x} - \mathbf{x}_{adv}\|_p \quad \text{s.t.} \quad f(\mathbf{x}) \neq f(\mathbf{x}_{adv})$$

## Optimal attack

$$\min_{\mathbf{x}_{adv}} \|\mathbf{x} - \mathbf{x}_{adv}\|_p \quad \text{s.t.} \quad f(\mathbf{x}) \neq f(\mathbf{x}_{adv})$$

## Kantchelian et al.

- proved the optimal attack on tree ensemble is NP-complete with increasing number of trees
- formulate the attack as a mixed integer linear program (MILP)

## Optimal attack

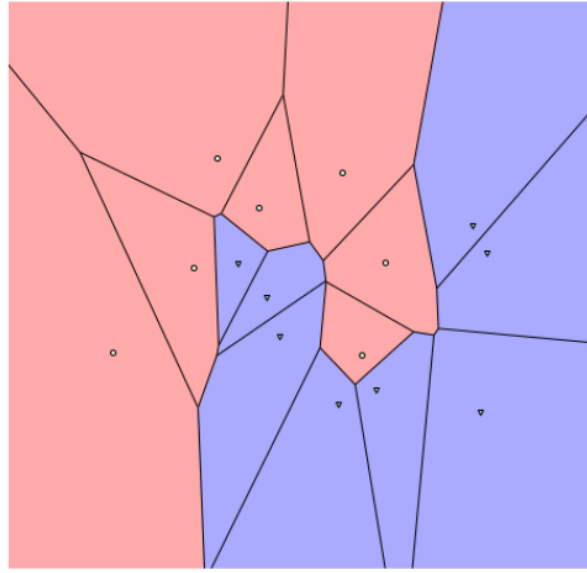
$$\min_{\mathbf{x}_{adv}} \|\mathbf{x} - \mathbf{x}_{adv}\|_p \quad \text{s.t.} \quad f(\mathbf{x}) \neq f(\mathbf{x}_{adv})$$

## Kantchelian et al.

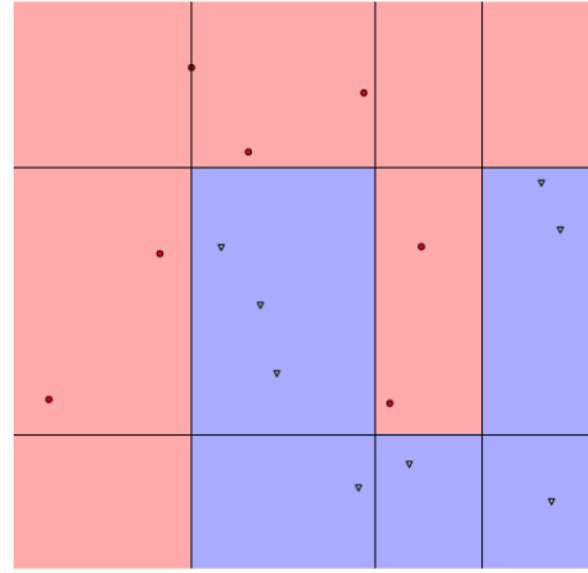
- proved the optimal attack on tree ensemble is NP-complete with increasing number of trees
- formulate the attack as a mixed integer linear program (MILP)

## Prior attacks on non-parametrics

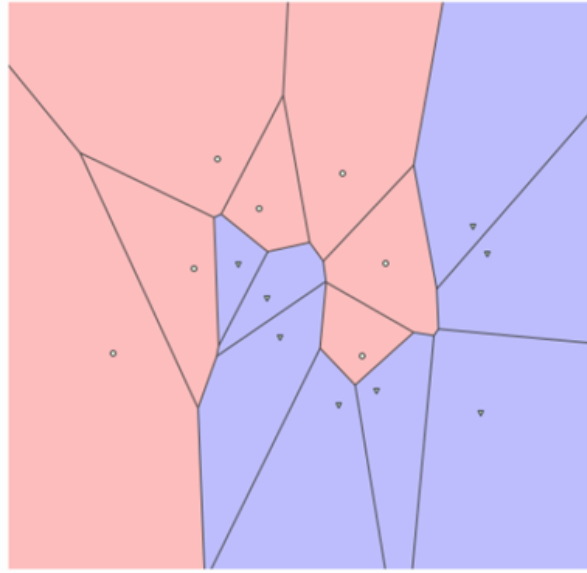
- model specific



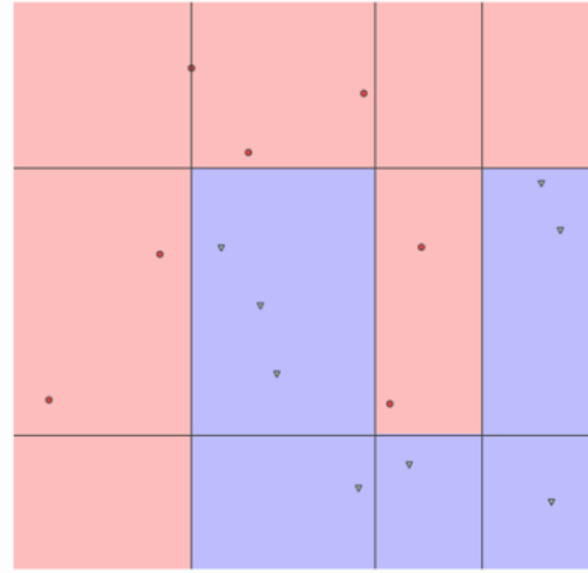
(a) 1-NN regions



(b) DT regions



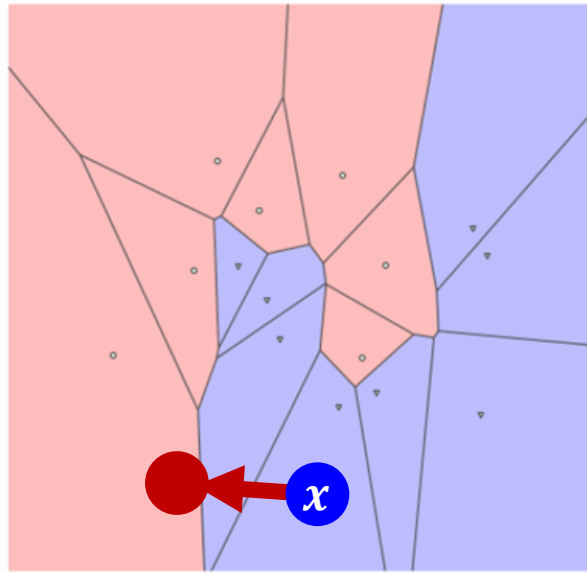
(a) 1-NN regions



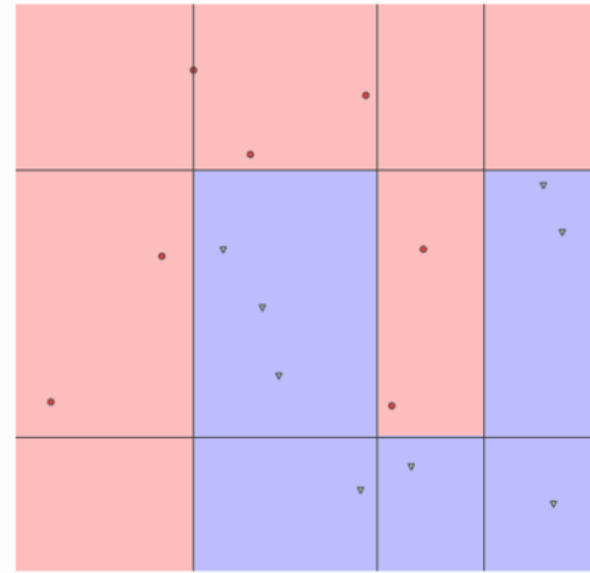
(b) DT regions

### Definition ( $(s, m)$ -decomposition)

The partition of  $\mathcal{R}^d$  into convex regions  $P_1, \dots, P_s$  s.t. each  $P_i$  can be described by at most  $m$  linear constraints.



(a) 1-NN regions



(b) DT regions

### Definition ( $(s, m)$ -decomposition)

The partition of  $\mathcal{R}^d$  into convex regions  $P_1, \dots, P_s$  s.t. each  $P_i$  can be described by at most  $m$  linear constraints.

## Region-Based Attack

$$\min_{i:f(\mathbf{x}) \neq y_i} \min_{\mathbf{x}_{adv} \in P_i} \|\mathbf{x} - \mathbf{x}_{adv}\|_p$$

- **outer min:** iterate through all regions
- **inner min:** LP for  $p = 1, \infty$  and QP for  $p = 2$

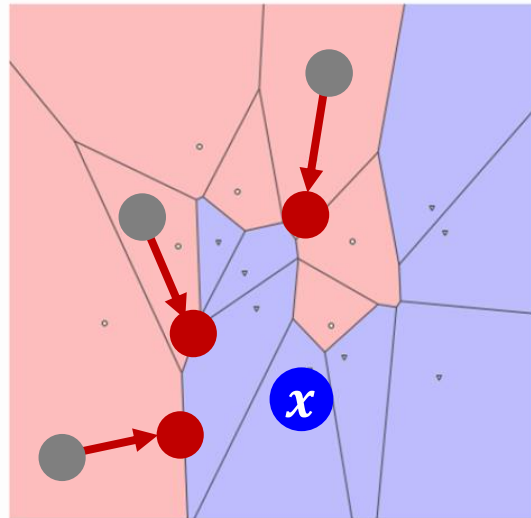
### Definition (*(s, m)-decomposition*)

The partition of  $\mathcal{R}^d$  into convex regions  $P_1, \dots, P_s$  s.t. each  $P_i$  can be described by at most  $m$  linear constraints.

# Region-Based Attack

$$\min_{i:f(\mathbf{x}) \neq y_i} \min_{\mathbf{x}_{adv} \in P_i} \|\mathbf{x} - \mathbf{x}_{adv}\|_p$$

- **outer min:** iterate through all regions
- **inner min:** LP for  $p = 1, \infty$  and QP for  $p = 2$



1-NN regions

## RBA-Exact:

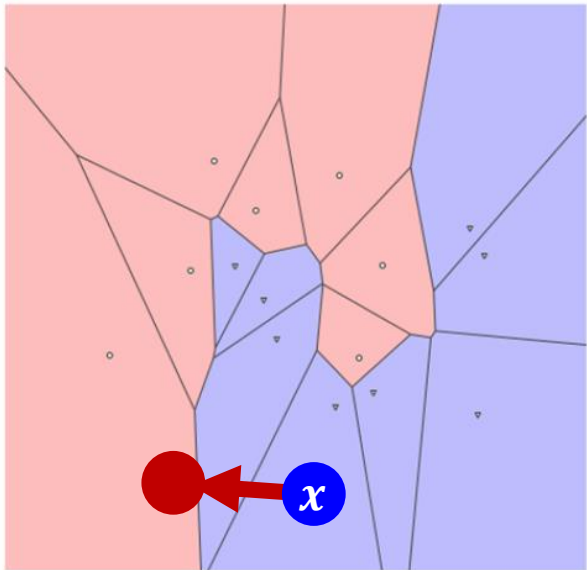
optimal adversarial example

## RBA-Approx:

only consider subset of 50-100 regions containing training pts



# Comparing Attacks



1-NN regions

Measure **distance to closest adversarial example** for examples in the test set

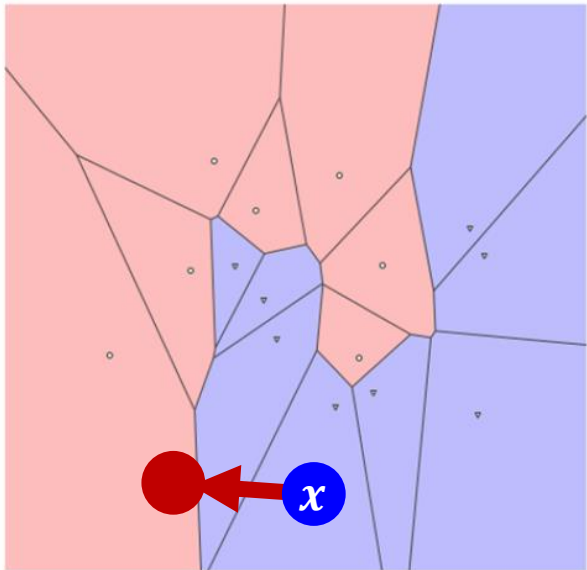
Lower == better attack

# Comparing Attacks

---

1-NN				3-NN			
Direct	BBox	Kernel	RBA-Exact	Direct	BBox	Kernel	RBA-Approx

---



1-NN regions

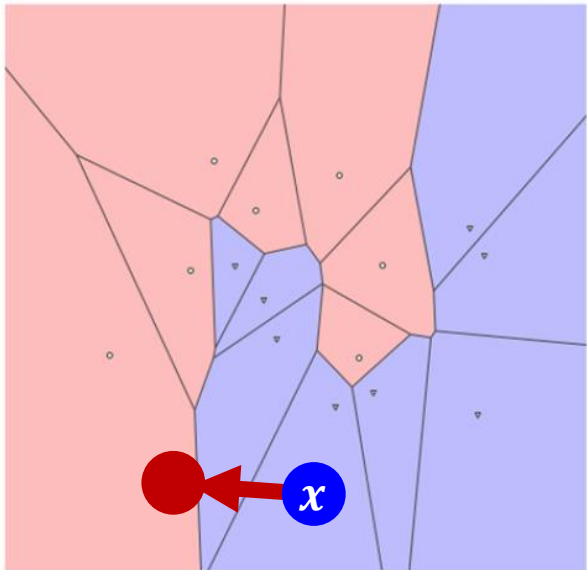
Measure **distance to closest adversarial example** for examples in the test set\*

\*after PCA to 25 dim. for MNIST/F-MNIST

Lower == better attack

# Comparing Attacks

1-NN			3-NN		
Direct	BBox	Kernel	Direct	BBox	Kernel
		<b>RBA-Exact</b>			<b>RBA-Approx</b>



1-NN regions

Measure **distance to closest adversarial example** for examples in the test set\*

\*after PCA to 25 dim. for MNIST/F-MNIST

Lower == better attack

# Comparing Attacks

	1-NN				3-NN			
	Direct	BBox	Kernel	RBA-Exact	Direct	BBox	Kernel	RBA-Approx
australian	.442	.336	.379	<b>.151</b>	.719	.391	.464	<b>.278</b>
cancer	.223	.364	.358	<b>.137</b>	.329	.376	.394	<b>.204</b>
covtype	.320	.207	.271	<b>.076</b>	.443	.265	.271	<b>.120</b>
diabetes	.074	.112	.165	<b>.035</b>	.130	.143	.191	<b>.078</b>
f-mnist06	.259	.162	.187	<b>.034</b>	.233	.184	.213	<b>.064</b>
f-mnist35	.354	.269	.288	<b>.089</b>	.355	.279	.295	<b>.111</b>
fourclass	.109	.124	.137	<b>.090</b>	.101	.113	.134	<b>.096</b>
halfmoon	.070	.129	.102	<b>.059</b>	.105	.132	.115	<b>.096</b>
mnist17	.330	.260	.239	<b>.079</b>	.302	.264	.247	<b>.098</b>

# Comparing Attacks

	1-NN				3-NN			
	Direct	BBox	Kernel	<b>RBA-Exact</b>	Direct	BBox	Kernel	<b>RBA-Approx</b>
australian	.442	.336	.379	<b>.151</b>	.719	.391	.464	<b>.278</b>
cancer	.223	.364	.358	<b>.137</b>	.329	.376	.394	<b>.204</b>
covtype	.320	.207	.271	<b>.076</b>	.443	.265	.271	<b>.120</b>
diabetes	.074	.112	.165	<b>.035</b>	.130	.143	.191	<b>.078</b>
f-mnist06	.259	.162	.187	<b>.034</b>	.233	.184	.213	<b>.064</b>
f-mnist35	.354	.269	.288	<b>.089</b>	.355	.279	.295	<b>.111</b>
fourclass	.109	.124	.137	<b>.090</b>	.101	.113	.134	<b>.096</b>
halfmoon	.070	.129	.102	<b>.059</b>	.105	.132	.115	<b>.096</b>
mnist17	.330	.260	.239	<b>.079</b>	.302	.264	.247	<b>.098</b>

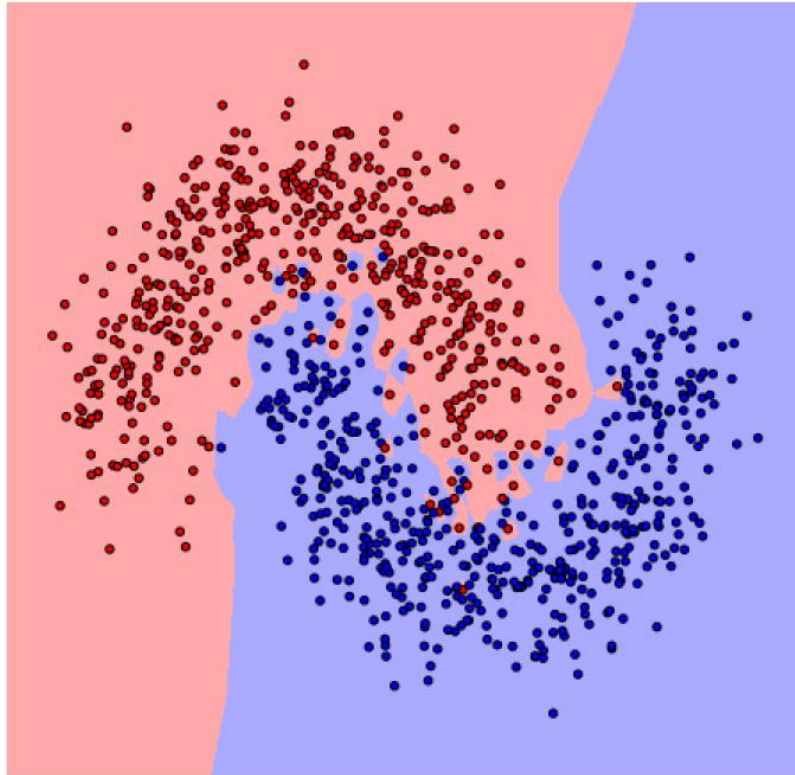
**Our attacks are 2-3x better**

# Comparing Attacks cont.

	Papernot's	DT BBox	RBA-Exact	BBox	RF RBA-Approx
australian	.140	.139	<b>.070</b>	<b>.364</b>	.446
cancer	.459	.334	<b>.255</b>	.451	<b>.383</b>
covtype	.289	.117	<b>.070</b>	.256	<b>.219</b>
diabetes	.237	.133	<b>.085</b>	<b>.181</b>	.184
f-mnist06	.200	.182	<b>.114</b>	.222	<b>.199</b>
f-mnist35	.287	.168	<b>.112</b>	<b>.201</b>	.246
fourclass	.288	.197	<b>.137</b>	.159	<b>.133</b>
halfmoon	.098	.148	<b>.085</b>	.182	<b>.149</b>
mnist17	.236	.175	<b>.117</b>	<b>.237</b>	.244

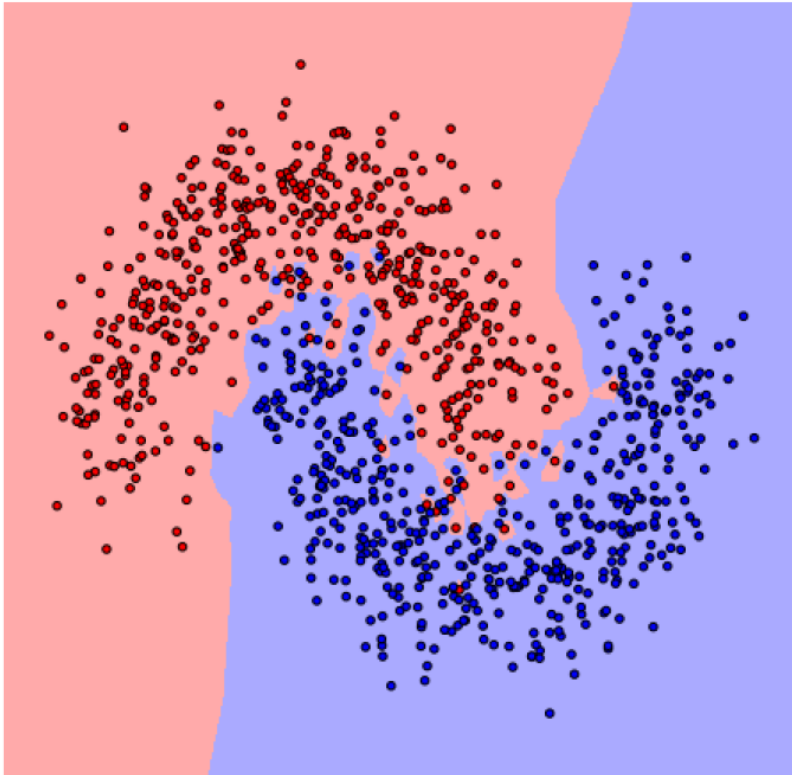
**RFs are much harder to attack**

# Improving robustness via separation

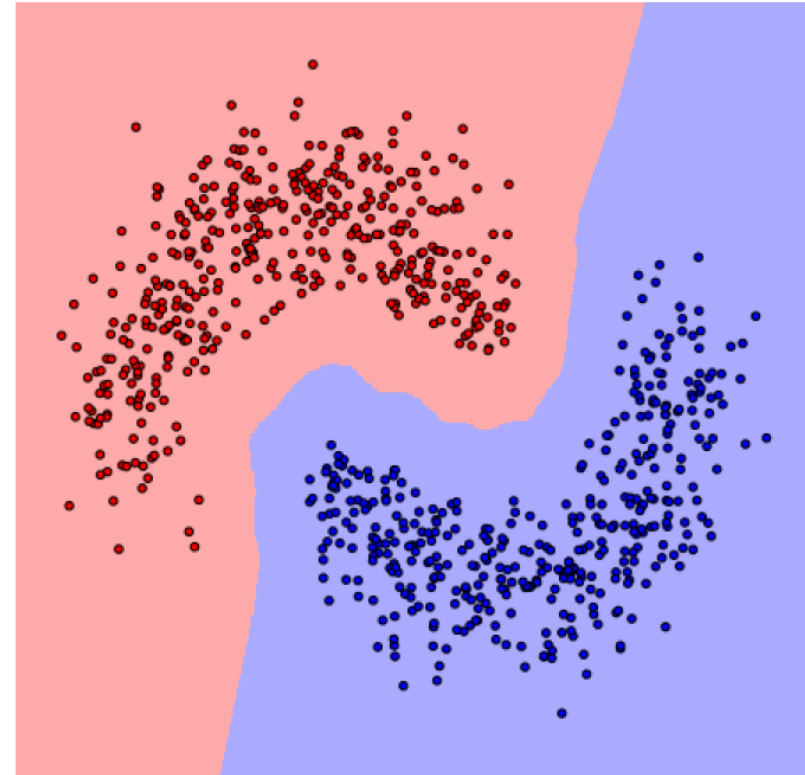


1-NN

# Improving robustness via separation



1-NN



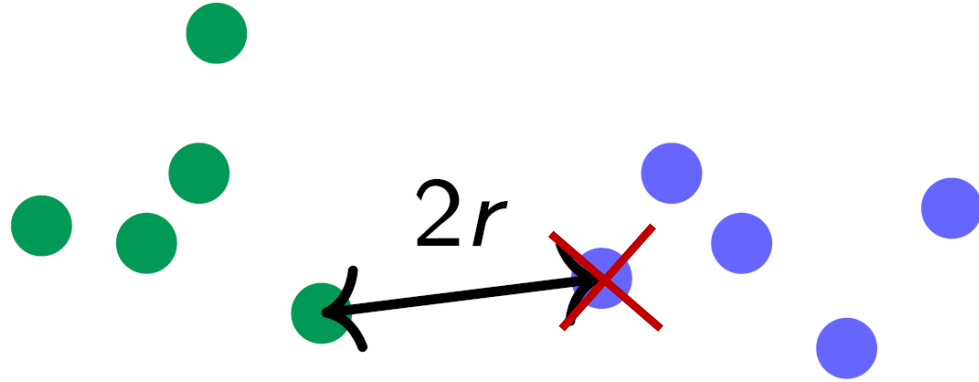
1-NN with separation (less overlap)



# Adversarial Pruning

## Defense strategy

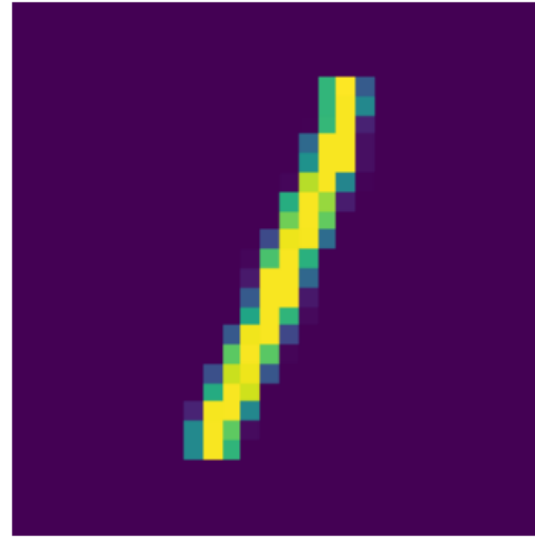
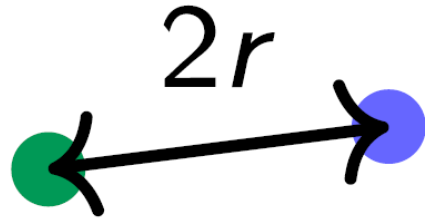
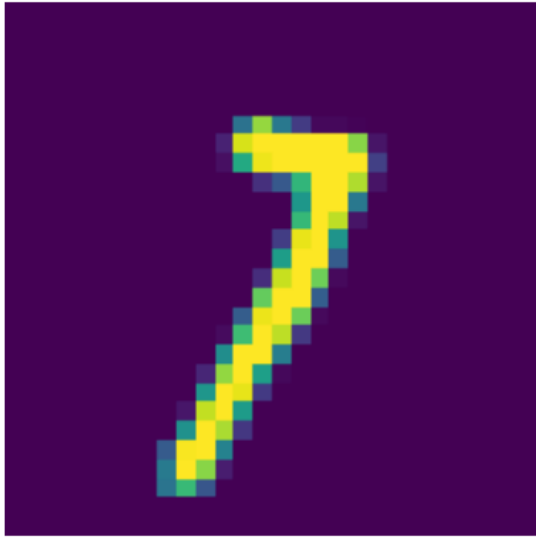
- 1 remove minimum # of examples s.t. distance between opposite labeled examples  $\geq 2r$
- 2 learn non-parametric classifier



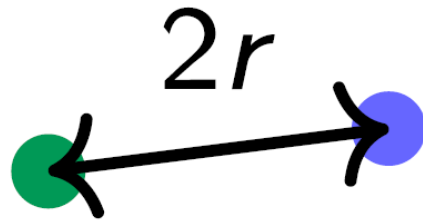
## Computing Pruned Dataset

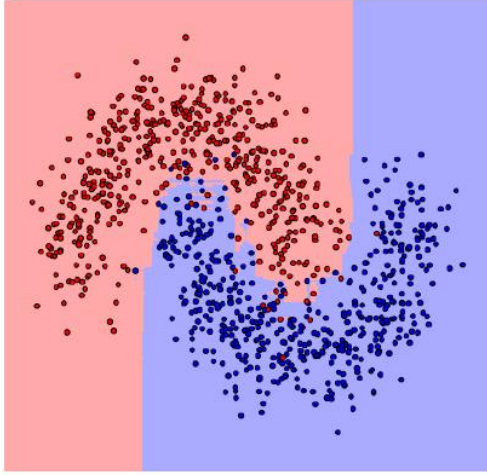
Bipartite maximum matching via Hopcroft-Karp algorithm (1973)

Graph  $n$  vertices and  $m$  edges, running time  $O(m \sqrt{n})$

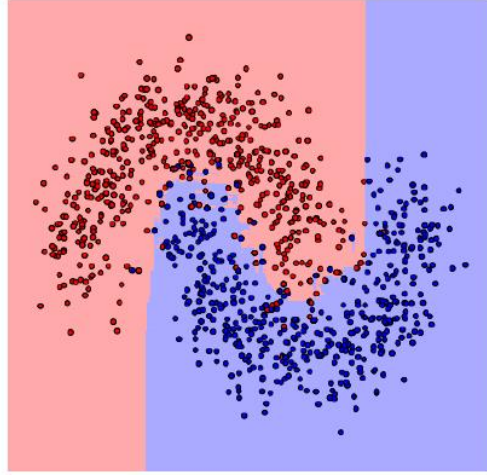


Distance in feature space after PCA to 25 dimensions

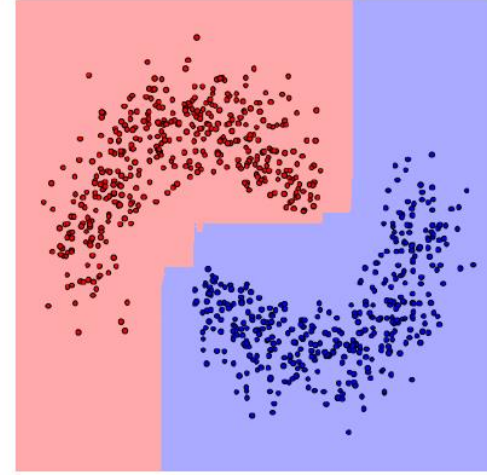




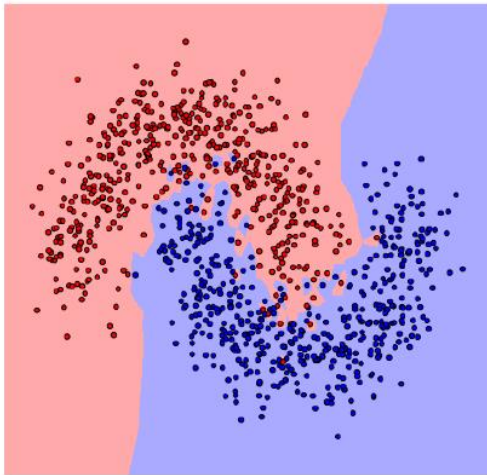
RF



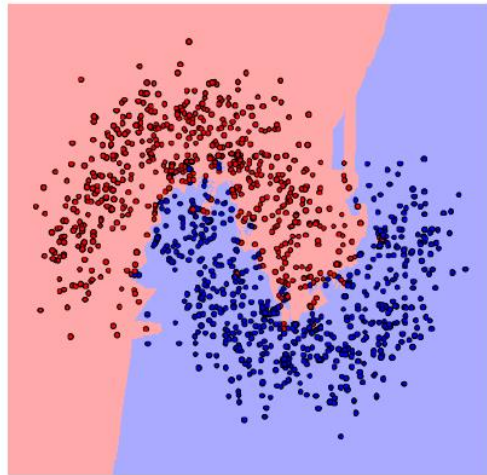
RF with AT



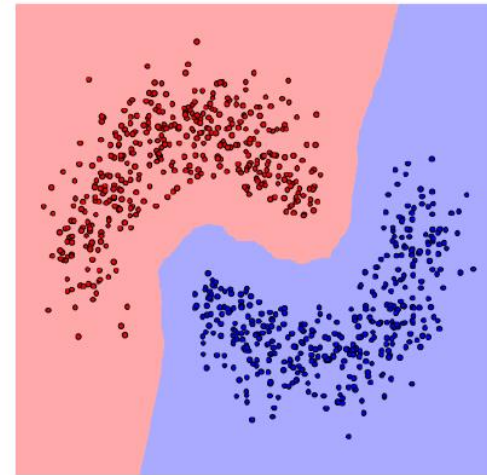
RF with AP



1-NN



1-NN with AT



1-NN with AP

# Evaluating defenses

---

		1-NN		3-NN		DT		RF			
	AT	Wang's	AP	AT	AP	AT	RS	AP	AT	RS	AP

---

---

Chen, Zhang, Boning, Hsieh. Robust Decision Trees Against Adversarial Examples. ICML 2018.

Wang, Jha, Chaudhuri. Analyzing the Robustness of Nearest Neighbors to Adversarial Examples. ICML 2018.

# Evaluating defenses

---

	AT	1-NN Wang's	AP		AT	3-NN AP		AT	DT RS	AP		AT	RF RS	AP
--	----	----------------	----	--	----	------------	--	----	----------	----	--	----	----------	----

---

$$\text{defscore} = \frac{\text{defended dist. to adv. example}}{\text{undefended dist. to adv. example}}$$

average over test set

restrict to correctly classified (normalize accuracy)

# Evaluating defenses

Baseline: Chen et al ICML'19  
RS = Robust splitting

---

		1-NN		3-NN		DT		RF		
	AT	Wang's	AP		AT	AP		AT	RS	AP

---

$$\text{defscore} = \frac{\text{defended dist. to adv. example}}{\text{undefended dist. to adv. example}}$$

average over test set

restrict to correctly classified (normalize accuracy)

# Evaluating defenses

Baseline: Chen et al ICML'19  
RS = Robust splitting

---

AT	1-NN Wang's	AP		AT	3-NN	AP		AT	DT RS	AP		AT	RF RS	AP
----	----------------	----	--	----	------	----	--	----	----------	----	--	----	----------	----

---

$$\text{defscore} = \frac{\text{defended dist. to adv. example}}{\text{undefended dist. to adv. example}}$$

average over test set

restrict to correctly classified (normalize accuracy)

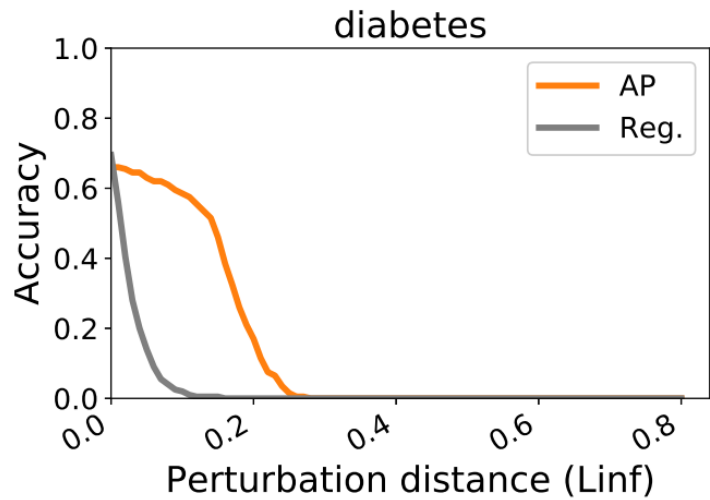
# Evaluating defenses

Baseline: Chen et al ICML'19  
RS = Robust splitting

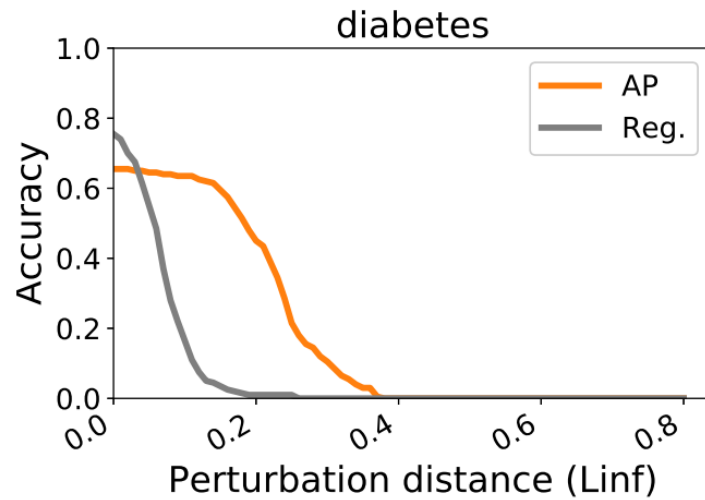
	AT	1-NN Wang's	AP	AT	3-NN AP	AT	DT RS	AP	AT	RF RS	AP
australian	0.64	<b>1.65</b>	<b>1.65</b>	0.68	<b>1.20</b>	2.36	<b>5.86</b>	2.37	1.07	<b>1.12</b>	1.04
cancer	0.82	1.05	<b>1.41</b>	1.06	<b>1.39</b>	0.85	1.09	<b>1.19</b>	0.87	<b>1.54</b>	1.26
covtype	0.61	<b>3.17</b>	<b>3.17</b>	0.81	<b>2.55</b>	1.07	2.90	<b>4.84</b>	0.93	1.59	<b>2.10</b>
diabetes	0.83	<b>4.69</b>	<b>4.69</b>	0.87	<b>2.97</b>	0.93	1.53	<b>2.22</b>	1.19	1.25	<b>2.22</b>
f-mnist06	0.94	2.09	<b>2.12</b>	0.86	<b>1.47</b>	0.82	<b>3.91</b>	1.85	0.97	1.17	<b>1.81</b>
f-mnist35	0.80	1.02	<b>1.08</b>	0.77	<b>1.05</b>	1.11	<b>2.64</b>	2.07	0.90	1.23	<b>1.32</b>
fourclass	0.93	<b>3.09</b>	<b>3.09</b>	0.89	<b>3.09</b>	1.06	1.23	<b>3.04</b>	1.03	1.92	<b>3.59</b>
halfmoon	1.03	1.98	<b>2.73</b>	0.93	<b>1.92</b>	1.54	1.98	<b>2.58</b>	1.04	1.01	<b>1.82</b>
mnist17	0.78	1.01	<b>1.20</b>	0.81	<b>1.13</b>	1.14	<b>2.91</b>	1.54	0.93	1.11	<b>1.29</b>

Higher == better defense



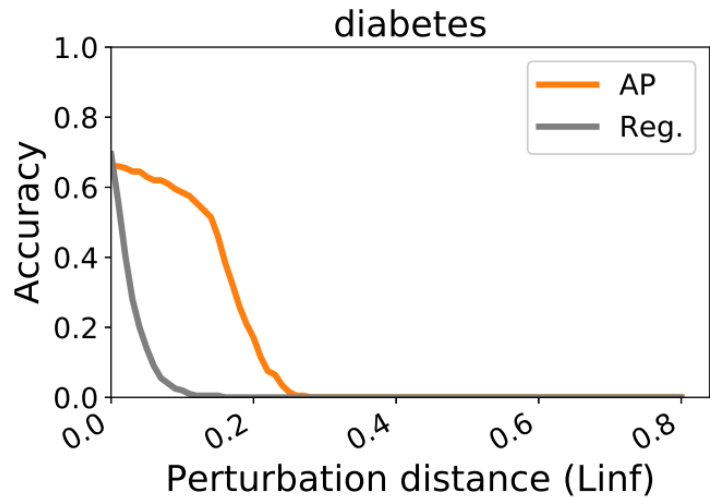


(e) 1-NN

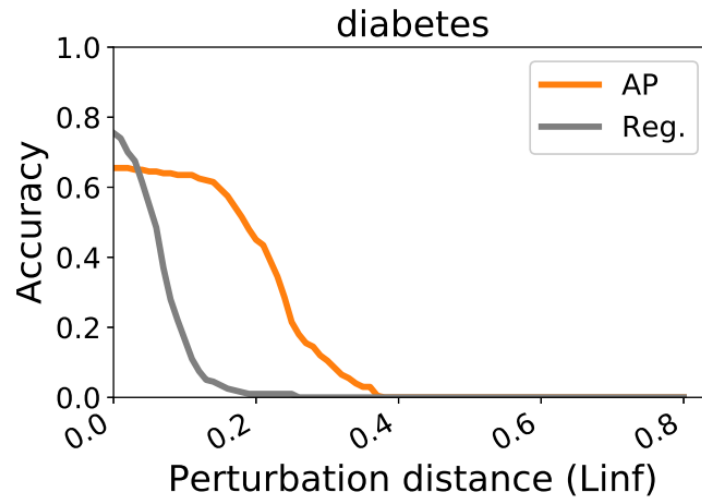


(f) 3-NN

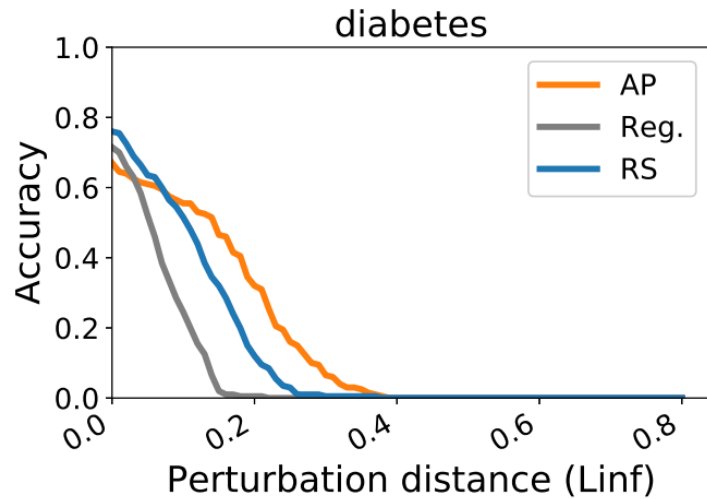
Our defense (AP)  
increases necessary  
perturbation distance



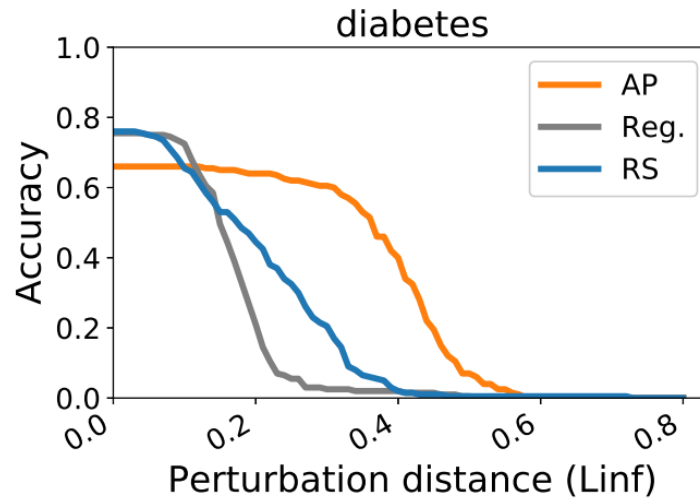
(e) 1-NN



(f) 3-NN



(g) Decision tree



(h) Random forest

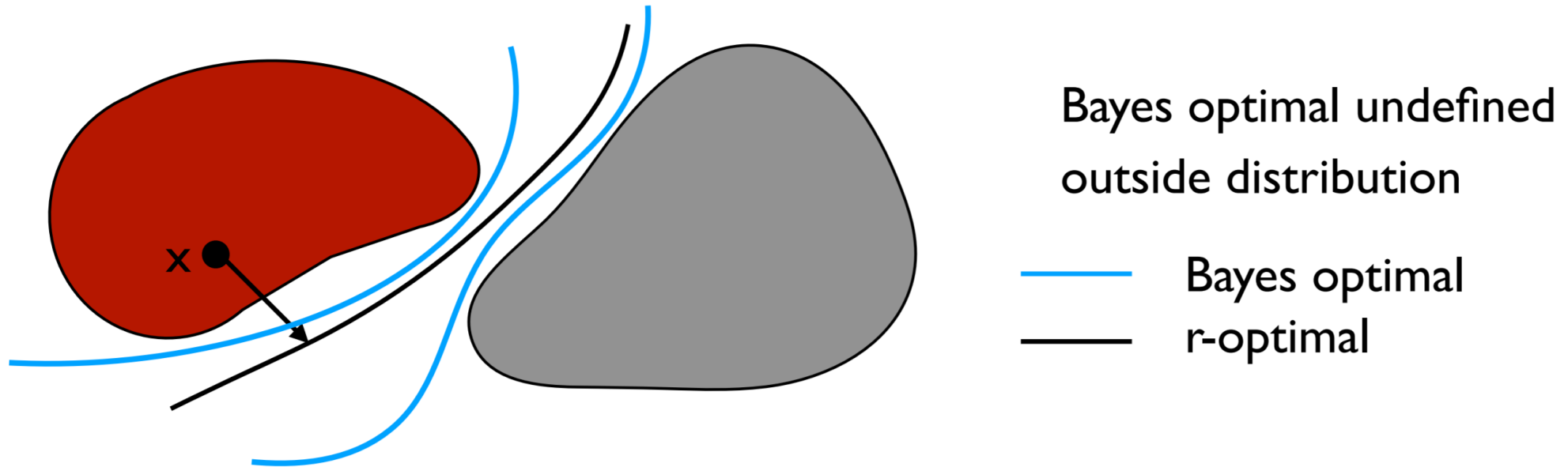
Our defense (AP) increases necessary perturbation distance

**Downside:**  
Reduces accuracy 😞

Theoretical justification:

Robust analogue of Bayes Optimal

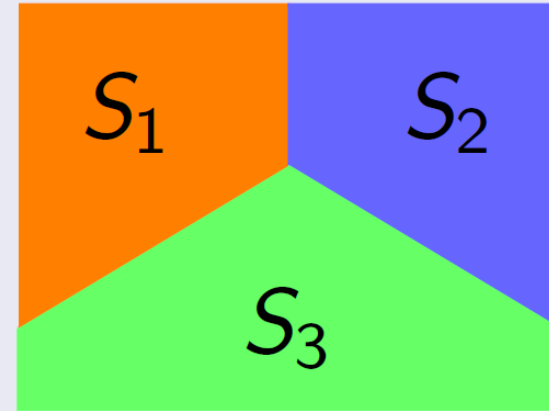
# Robust analogue of Bayes Optimal



**$r$ -optimal** = classifier that maximizes accuracy at points that have robustness radius at least  $r$

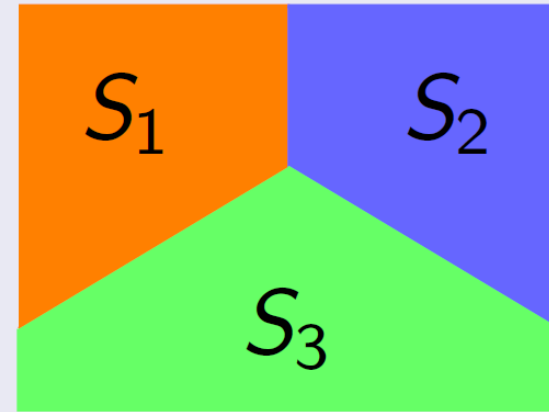
# Bayes-optimal classifier

$$\max_{S_1, \dots, S_c} \sum_{j=1}^c \int_{\mathbf{x} \in S_j} pr(y = j | \mathbf{x}) d\mu$$



## Bayes-optimal classifier

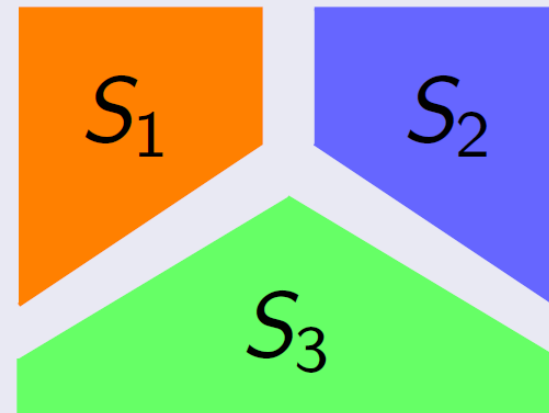
$$\max_{S_1, \dots, S_c} \sum_{j=1}^c \int_{\mathbf{x} \in S_j} pr(y = j | \mathbf{x}) d\mu$$



## $r$ -Optimal classifier

$$\max_{S_1, \dots, S_c} \sum_{j=1}^c \int_{\mathbf{x} \in S_j} pr(y = j | \mathbf{x}) d\mu$$

$$\text{s.t. } d(S_j, S_{j'}) \geq 2r \quad \forall j \neq j'$$



## Attack algorithm

- target model  $f$
- target example  $\mathbf{x}$
- attack budget  $r$  (defines “small”)

attack algorithm  $A(f, \mathbf{x}, r) \rightarrow \mathcal{R}^d$  returns an example under some attack budget constraint  $r$

## Definition (Astuteness $\text{ast}_\mu(A, f, r)$ )

The accuracy after attack. Let  $\mu$  be a distribution on  $\mathcal{R}^d \times c$  and

$$\text{ast}_\mu(A, f, r) := \Pr_{(\mathbf{x}, y) \sim \mu} [f(\mathbf{x}) = f(A(f, \mathbf{x}, r)) \text{ and } f(\mathbf{x}) = y]$$

## Attack algorithm

- target model  $f$
- target example  $\mathbf{x}$
- attack budget  $r$  (defines “small”)

attack algorithm  $A(f, \mathbf{x}, r) \rightarrow \mathcal{R}^d$  returns an example under some attack budget constraint  $r$

## Definition (Astuteness $\text{ast}_\mu(A, f, r)$ )

The accuracy after attack. Let  $\mu$  be a distribution on  $\mathcal{R}^d \times c$  and

$$\text{ast}_\mu(A, f, r) := \Pr_{(\mathbf{x}, y) \sim \mu} [f(\mathbf{x}) = f(A(f, \mathbf{x}, r)) \text{ and } f(\mathbf{x}) = y]$$

## Theorem

*$r$ -Optimal classifier maximizes astuteness with attack radius  $r$  under  $\mu$ .*

$$f_{\text{ropt}} = \operatorname{argmax}_f \text{ast}_\mu(f, r)$$



# Follow-up Theory [Bhattacharjee, Chaudhuri 2020]

They prove that **Adversarial Pruning** +  $k$ -NN or kernel classifiers converges toward **optimally robust** and accurate classifiers (under certain conditions)

## Definition (Astuteness $\text{ast}_\mu(A, f, r)$ )

The accuracy after attack. Let  $\mu$  be a distribution on  $\mathcal{R}^d \times c$  and

$$\text{ast}_\mu(A, f, r) := \Pr_{(\mathbf{x}, y) \sim \mu} [f(\mathbf{x}) = f(A(f, \mathbf{x}, r)) \text{ and } f(\mathbf{x}) = y]$$

## Theorem

*$r$ -Optimal classifier maximizes astuteness with attack radius  $r$  under  $\mu$ .*

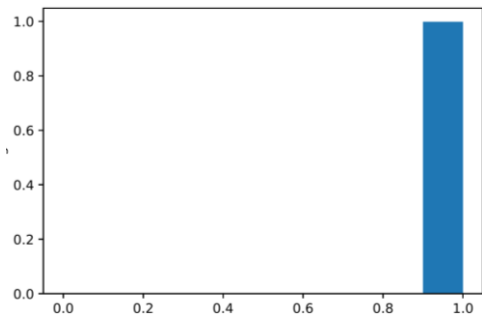
$$f_{\text{ropt}} = \operatorname{argmax}_f \text{ast}_\mu(f, r)$$

What about neural networks?

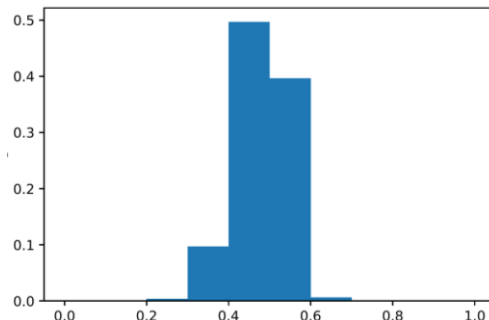
Can we get robustness + accuracy?

# Separation of real datasets

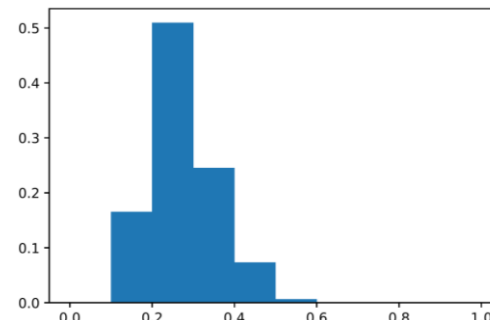
## MNIST



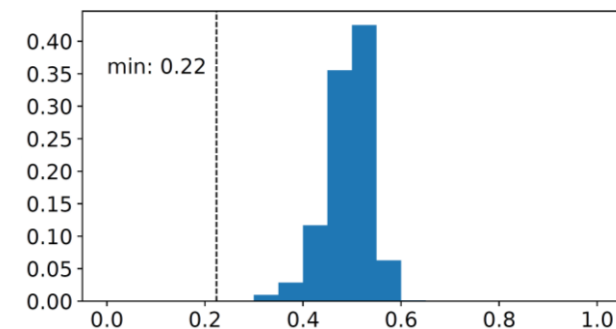
## CIFAR-10



## SVHN



## ResImageNet



pairwise  $L_\infty$  distance

---

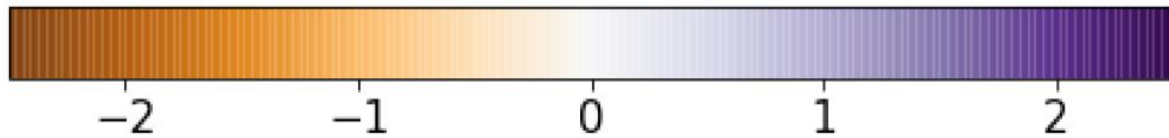
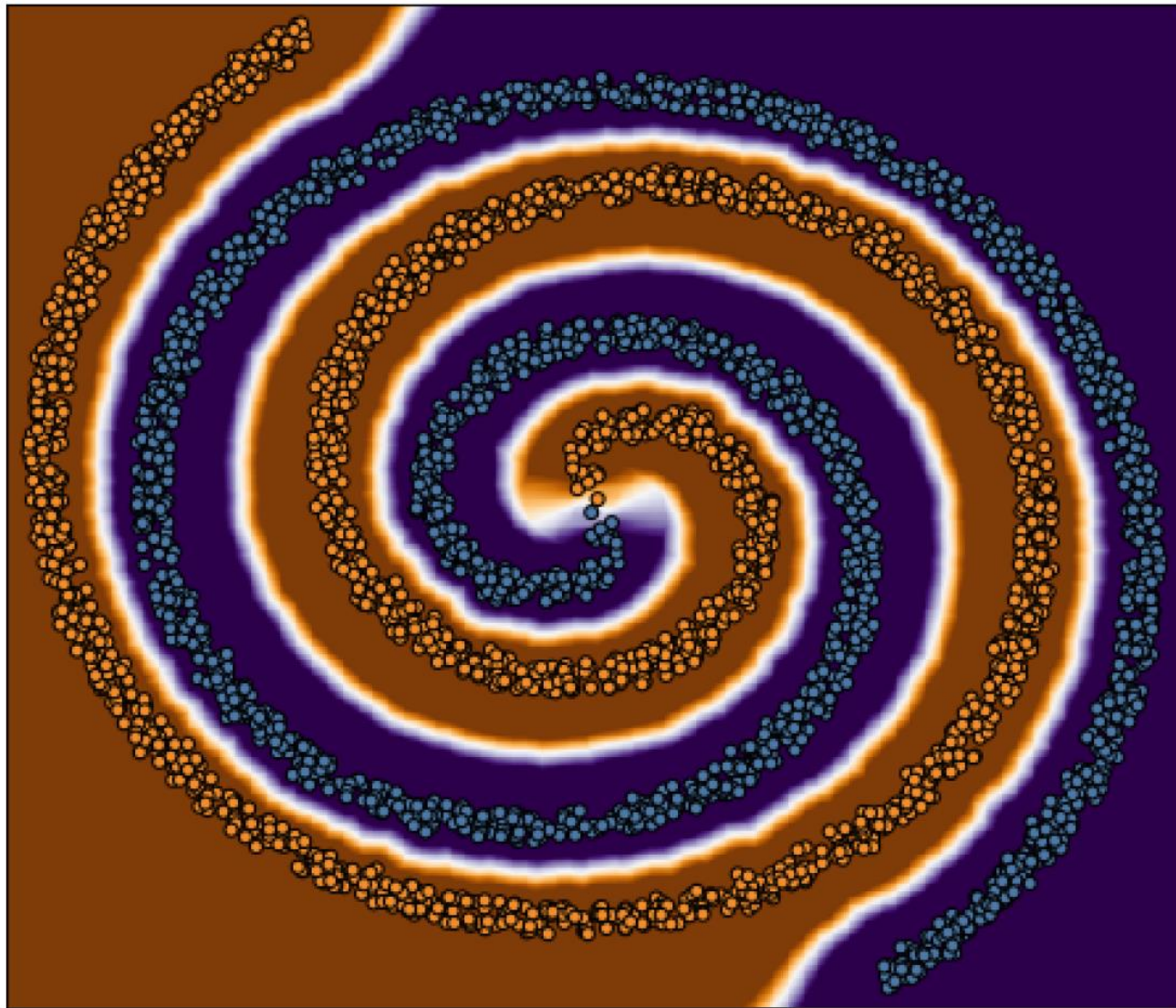
	typical perturbation distance	minimum Train-Train separation	minimum Test-Train separation
MNIST	0.1	0.737	0.812
CIFAR-10	0.031 (8/255)	0.212	0.220
SVHN	0.031 (8/255)	0.094	0.110
Restricted ImageNet	0.005	0.180	0.224

---

## **Thm.**

For separated data there always exists a classifier that is

- Accurate
- Robust
- Locally Lipschitz



## Thm.

For separated data there always exists a classifier that is

- Accurate
- Robust
- Locally Lipschitz

# Perfect accuracy & robustness, at least in theory

## Locally Lipschitz:

A function  $f$  is  $L$ -Locally Lipschitz in a radius  $r$  around  $x$  if for all  $x'$  s.t.  $d(x, x') \leq r$ , we have  $|f(x) - f(x')| < L \cdot d(x, x')$

## Key idea behind all (provable) results for adversarial robustness

---

Hein, Andriushchenko. Formal Guarantees on the Robustness of a Classifier Against Adversarial Manipulation. NeurIPS 2017.

Cohen, Rosenfeld, Kolter. Certified Adversarial Robustness via Randomized Smoothing. ICML 2019.

Salman, Yang, Li, Zhang, Zhang, Razenshteyn, Bubeck. Provably Robust Deep Learning via Adversarially Trained Smoothed Classifiers. NeurIPS 2019.

# Perfect accuracy & robustness, at least in theory

## Locally Lipschitz:

A function  $f$  is  $L$ -Locally Lipschitz in a radius  $r$  around  $x$  if for all  $x'$  s.t.  $d(x, x') \leq r$ , we have  $|f(x) - f(x')| < L \cdot d(x, x')$

## Key idea behind all (provable) results for adversarial robustness

**Lemma.** If  $f$  is  $L$ -Locally Lipschitz, then  $g = \text{sign}(f)$  is robust at  $x$  whenever  $|f(x)| \geq \frac{1}{Lr}$

---

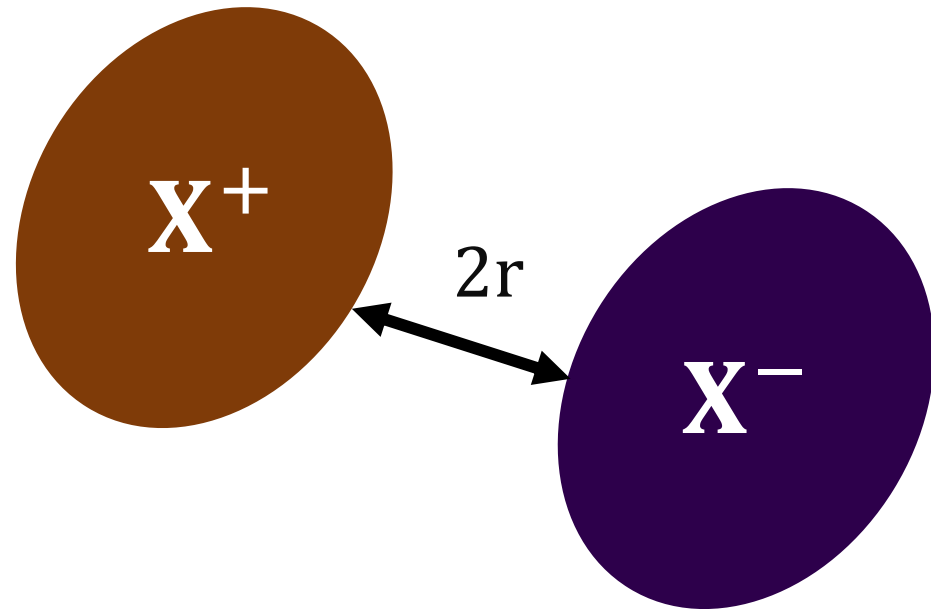
Hein, Andriushchenko. Formal Guarantees on the Robustness of a Classifier Against Adversarial Manipulation. NeurIPS 2017.

Cohen, Rosenfeld, Kolter. Certified Adversarial Robustness via Randomized Smoothing. ICML 2019.

Salman, Yang, Li, Zhang, Zhang, Razenshteyn, Bubeck. Provably Robust Deep Learning via Adversarially Trained Smoothed Classifiers. NeurIPS 2019.

# Perfect accuracy & robustness, at least in theory

**Theorem.** If data is  $2r$  – separated, there always exists a classifier that is perfectly robust and accurate, based on a function with local Lipschitz constant  $1/r$ .





# Perfect accuracy & robustness, at least in theory

**Theorem.** If data is  $2r$  – separated, there always exists a classifier that is perfectly robust and accurate, based on a function with local Lipschitz constant  $1/r$ .

**Proof.** Classifier:  $g = \text{sign}(f)$  where  $f(\mathbf{x}) = \frac{d(\mathbf{x}, \mathcal{X}^-) - d(\mathbf{x}, \mathcal{X}^+)}{2r}$

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{x}') &= \frac{d(\mathbf{x}, \mathcal{X}^-) - d(\mathbf{x}', \mathcal{X}^-) - d(\mathbf{x}, \mathcal{X}^+) + d(\mathbf{x}', \mathcal{X}^+)}{2r} \\ &\leq \frac{2d(\mathbf{x}, \mathbf{x}')}{2r} \end{aligned}$$

**Lemma**  $\rightarrow$  Robust + Accurate if  $|f(\mathbf{x})| \geq 1$

# Robustness is more convoluted...

**Adversarial Training:**  $\min_f \mathbb{E} \left\{ \max_{\mathbf{X}' \in \mathbb{B}(\mathbf{X}, \epsilon)} \mathcal{L}(f(\mathbf{X}'), Y) \right\}$

**Gradient Regularization:**  $\min_f \mathbb{E} \left\{ \mathcal{L}(f(\mathbf{X}), Y) + \beta \|\nabla_{\mathbf{X}} \mathcal{L}(f(\mathbf{X}), Y)\|_2^2 \right\}$

**TRADES:**  $\min_f \mathbb{E} \left\{ \mathcal{L}(f(\mathbf{X}), Y) + \beta \max_{\mathbf{X}' \in \mathbb{B}(\mathbf{X}, \epsilon)} \mathcal{L}(f(\mathbf{X}), f(\mathbf{X}')) \right\}$

---

Madry, Makelov, Schmidt, Tsipras, Vladu. Towards Deep Learning Models Resistant to Adversarial Attacks. ICML 2018.

Finlay, Oberman. Scalable Input Gradient Regularization for Adversarial Robustness. 2019.

Zhang, Yu, Jiao, Xing, Ghaoui, Jordan. Theoretically Principled Trade-off Between Robustness and Accuracy. ICML 2019.

# CIFAR-10 Results

methods	train accuracy	test accuracy	adv test accuracy	test lipschitz
Natural	100.00	93.81	0.00	425.71
GR	94.90	80.74	21.32	28.53
LLR	100.00	91.44	22.05	94.68
AT	99.84	83.51	43.51	26.23
TRADES( $\beta=1$ )	99.76	84.96	43.66	28.01
TRADES( $\beta=3$ )	99.78	85.55	46.63	22.42
TRADES( $\beta=6$ )	98.93	84.46	48.58	13.05

# Conclusion

New attacks + defense for  $k$ -NN, DT, RF

Theory for well-separated data and local Lipschitzness

**Q1:** Better attack/defense for Random Forests?

**Q2:** How to achieve high accuracy AND robustness?

**Q3:** Beyond local Lipschitzness for provable robustness?

# Thanks!

Cyrus Rashtchian

[www.cyrusrashtchian.com](http://www.cyrusrashtchian.com)

UCSD

new blog: [ucsdml.github.io](http://ucsdml.github.io)



@CyrusRashtchian

kitten or ice cream





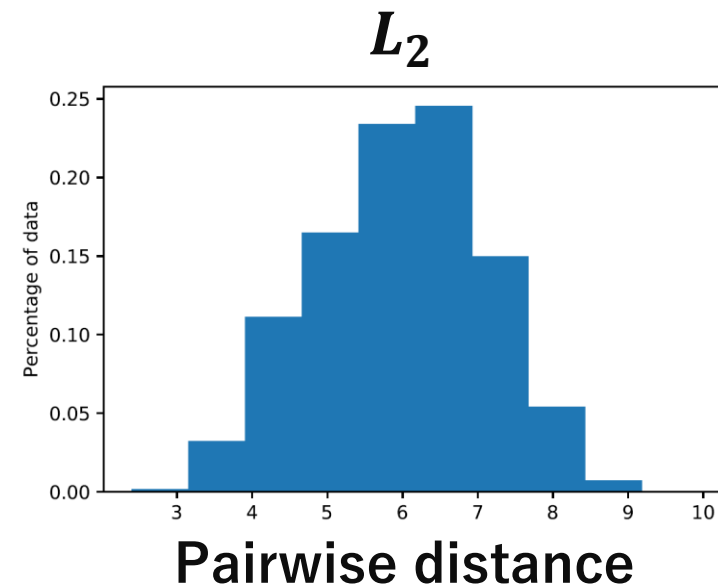
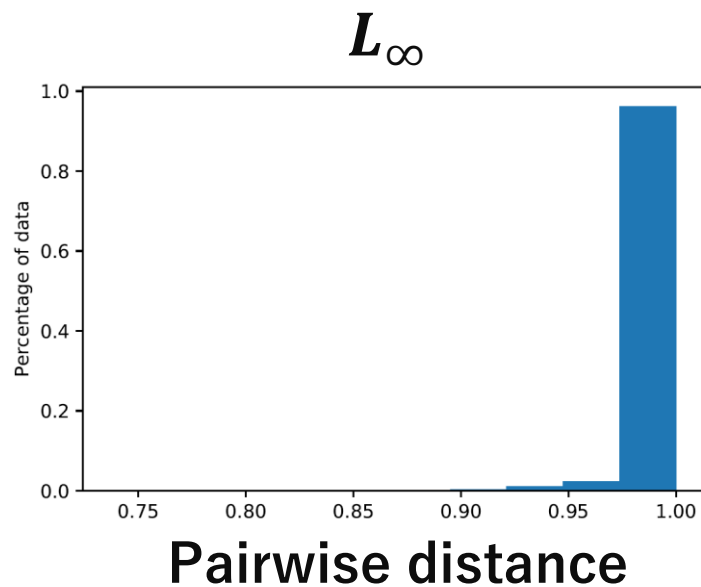
# Restricted ImageNet Dataset (1.3M Images)

Classes used in the Restricted ImageNet model. The class ranges are inclusive.

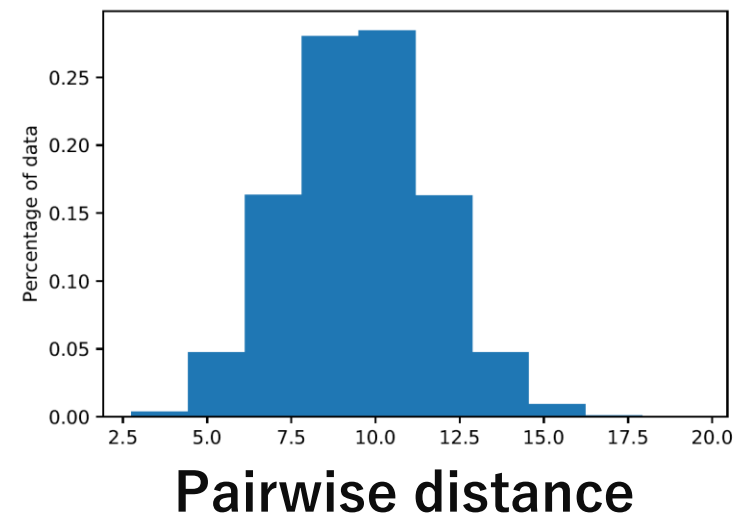
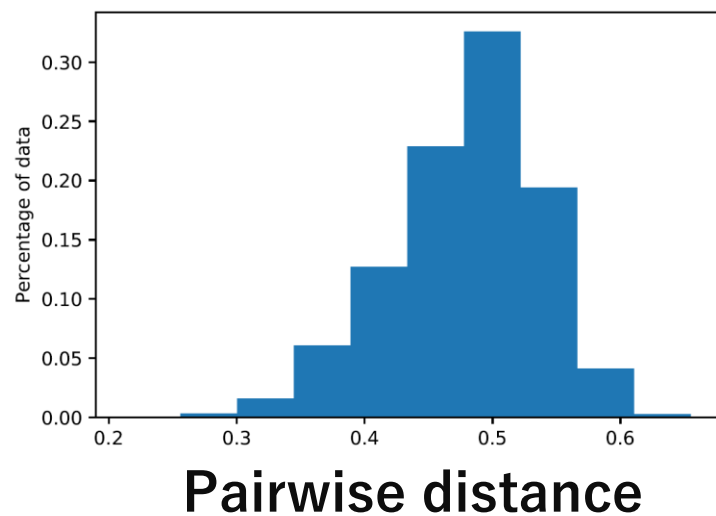
<b>Class</b>	<b>Corresponding ImageNet Classes</b>
"Dog"	151 to 268
"Cat"	281 to 285
"Frog"	30 to 32
"Turtle"	33 to 37
"Bird"	80 to 100
"Primate"	365 to 382
"Fish"	389 to 397
"Crab"	118 to 121
"Insect"	300 to 319

# Separation of real datasets

**MNIST**



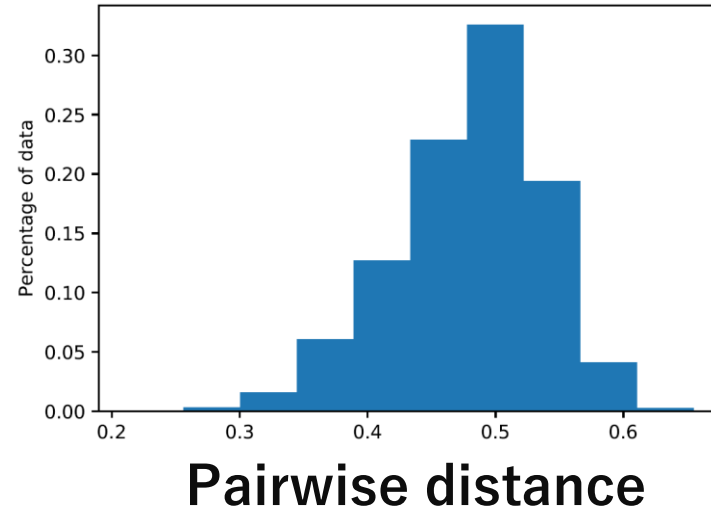
**CIFAR-10**



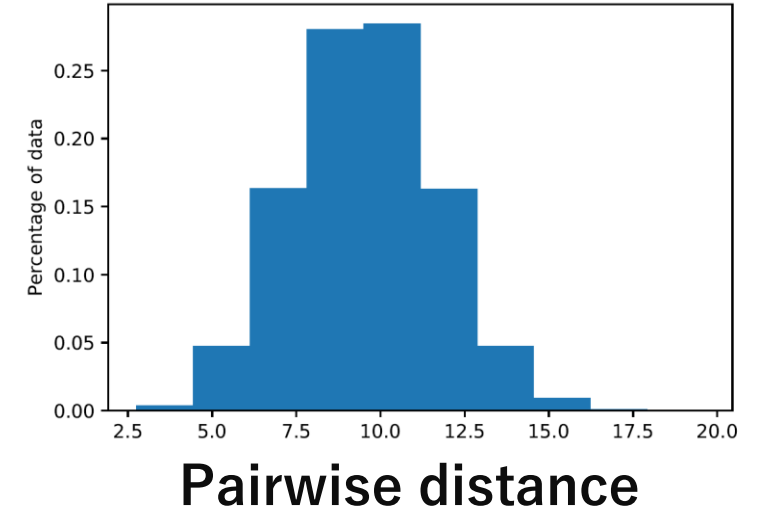


# CIFAR-10

$L_\infty$



$L_2$



# MNIST Results

methods	train accuracy	test accuracy	adv test accuracy	test lipschitz
Natural	100.00	99.20	59.83	67.25
GR	99.99	99.29	91.03	26.05
LLR	100.00	99.43	92.14	30.44
AT	99.98	99.31	97.21	8.84
TRADES( $\beta=1$ )	99.81	99.26	96.60	9.69
TRADES( $\beta=3$ )	99.21	98.96	96.66	7.83
TRADES( $\beta=6$ )	97.50	97.54	93.68	2.87

# ResImageNet Results

Restricted ImageNet	train accuracy	test accuracy	adv test accuracy	test lipschitz
Natural	97.72	93.47	7.89	32228.51
GR	91.12	88.51	62.14	886.75
LLR	98.76	93.44	52.65	4795.66
AT	96.22	90.33	82.25	287.97
TRADES( $\beta=1$ )	97.39	92.27	79.90	2144.66
TRADES( $\beta=3$ )	95.74	90.75	82.28	396.67
TRADES( $\beta=6$ )	93.34	88.92	82.13	200.90

# NP Hard for RFs

To find any adversarial example

Reduction from 3SAT

# clauses = # trees

(depth 3)

$$(x_0 \vee \neg x_1 \vee x_2) \wedge (x_1 \vee x_3 \vee \neg x_4) \wedge \dots$$

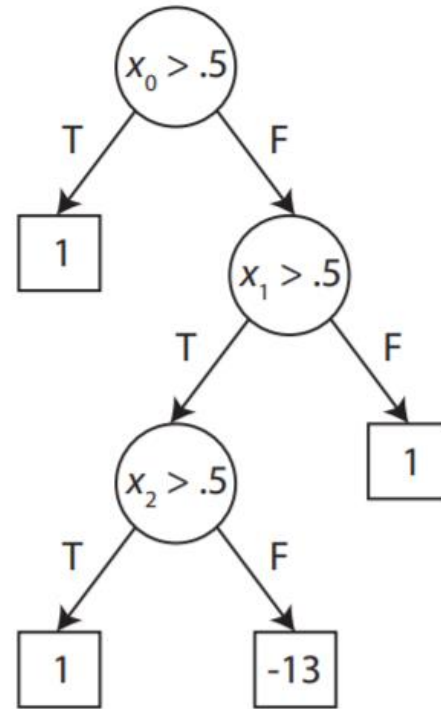
# NP Hard for RFs

To find any adversarial example

Reduction from 3SAT

# clauses = # trees

(depth 3)



$$(x_0 \vee \neg x_1 \vee x_2)$$

**+1 if clause is True**  
**-1 if clause is False**

$$(x_0 \vee \neg x_1 \vee x_2) \wedge (x_1 \vee x_3 \vee \neg x_4) \wedge \dots$$

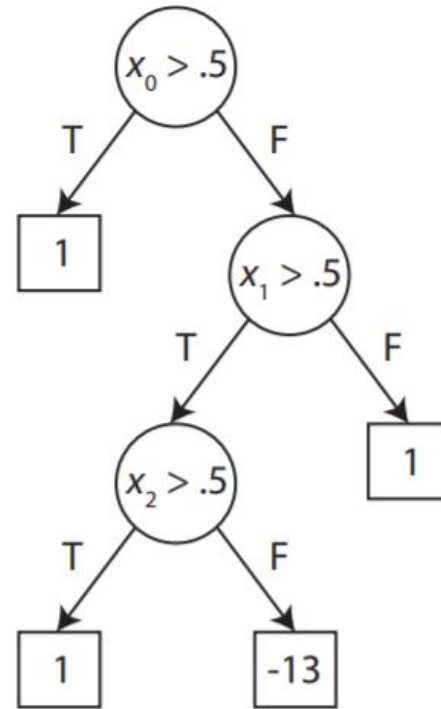
# NP Hard for RFs

To find any adversarial example

Reduction from 3SAT

# clauses = # trees

(depth 3)



**RF outputs True**

if and only if

**formula is SAT**

$$(x_0 \vee \neg x_1 \vee x_2)$$

**+1 if clause is True**  
**-T if clause is False**

$$(x_0 \vee \neg x_1 \vee x_2) \wedge (x_1 \vee x_3 \vee \neg x_4) \wedge \dots$$