Limits on Robustness to Adversarial Examples

Elvis Dohmatob

Criteo Al Lab

October 2, 2019



Table of contents



2 Classifier-dependent lower bounds



Preliminaries on adversarial robustness

Definition of adversarial attacks

A classifier is trained and deployed (e.g the computer vision system on a self-driving car)

 $\blacksquare At$ test / inference time, an attacker may submit queries to the classifier by

- sampling a real sample point x with true label k (e.g "pig"),
- modifying it $x \mapsto x^{adv}$ given to a prescribed threat model.

Goal of attacker is to make classifier label x^{adv} as $\neq k$ (e.g airliner)

The flying pig!



(Picture is courtesy of https://gradientscience.org/intro_adversarial/)

■ $x \mapsto x^{adv} := x + \text{ noise}$, $\|\text{noise}\| \le \varepsilon = 0.005$ (in example above) ■ Fast Gradient Sign Method: noise = sign($\nabla_x loss(h(x), y)$)

FGSM for generating adversarial examples [Goodfellow '14]

109 • clas	ss WI	<pre>iteBoxAttack(object):</pre>
		; adversarial perturbation of features (in sup norm threat model) a first-order Taylor approx. of the loss function, w.r.t the features
	ает	<pre>init(setf, epsilon): self.epsilon = epsilon</pre>
	def	<pre>call(self, model, true_features, true_labels, loss_func=F.nll_loss):</pre>
		true_reatures = variable(true_reatures, requires_grad=irue)
		pred_labels = model.torward(true_teatures)
		loss = loss_func(pred_labels, true_labels)
		<pre>grad = autograd.grad(loss, true features, retain graph=True)[0]</pre>
		<pre>return true_features + self.epsilon * torch.sign(grad) # move uphill</pre>

 $\blacksquare x \mapsto x^{\mathsf{adv}} := \mathsf{clip}(x + \varepsilon \mathsf{sign}(\nabla_x \mathsf{loss}(h(x), y)))$

Adversarial attacks and defenses, an arms race!



Image courtesy of [Goldstein' 19; Shafahi '19]

Elvis Dohmatob

Limits on Robustness to Adversarial Examples - slide 7 / 41

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Classifier-dependent lower bounds

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Problem setup

A classifier is simply a Borel-measurable mapping $h : \mathcal{X} \to \mathcal{Y}$ from feature space \mathcal{X} (with metric *d*) to label space $\mathcal{Y} := \{1, \dots, K\}$.

A classifier is trained and deployed (e.g the computer vision system on a self-driving car)

At test / inference time, an attacker may submit queries to the classifier by sampling a real sample point $x \in \mathcal{X}$ with true label $k \in \mathcal{Y}$, and modifying it $x \mapsto x^{adv}$ according to a prescribed threat model.

- For example, modifying a few pixels on a road traffic sign [Su et al. '17]
- Modifying intensity of pixels by a limited amount determined by a prescribed tolerance level [Tsipras '18], etc., on it.

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Problem setup

A classifier is simply a Borel-measurable mapping $h : \mathcal{X} \to \mathcal{Y}$ from feature space \mathcal{X} (with metric *d*) to label space $\mathcal{Y} := \{1, \dots, K\}$.

A classifier is trained and deployed (e.g the computer vision system on a self-driving car)

At test / inference time, an attacker may submit queries to the classifier by sampling a real sample point $x \in \mathcal{X}$ with true label $k \in \mathcal{Y}$, and modifying it $x \mapsto x^{adv}$ according to a prescribed threat model.

- For example, modifying a few pixels on a road traffic sign [Su et al. '17]
- Modifying intensity of pixels by a limited amount determined by a prescribed tolerance level [Tsipras '18], etc., on it.

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Problem setup: notations

Standard accuracy: $\operatorname{acc}(h|k) := 1 - \operatorname{err}(h|k)$, where $\operatorname{err}(h|k) := P_{X|k}(h(X) \neq k)$ is the error of *h* on class *k*.

• Small $\operatorname{acc}(h|k) \implies h$ is inaccurate on class k.

■ Adversarial robustness accuracy: $\operatorname{acc}_{\varepsilon}(h|k) := 1 - \operatorname{err}_{\varepsilon}(h|k)$, where $\operatorname{err}_{\varepsilon}(h|k) := P_{X|k}(\exists x' \in \operatorname{Ball}(X; \varepsilon) \mid h(x') \neq k)$ is the adversarial robustness error of h on class k.

• Small $\operatorname{acc}_{\varepsilon}(h|k) \implies h$ is vulnerable to attacks on class k.

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Problem setup: notations

Standard accuracy: $\operatorname{acc}(h|k) := 1 - \operatorname{err}(h|k)$, where $\operatorname{err}(h|k) := P_{X|k}(h(X) \neq k)$ is the error of *h* on class *k*.

• Small $\operatorname{acc}(h|k) \implies h$ is inaccurate on class k.

■ Adversarial robustness accuracy: $\operatorname{acc}_{\varepsilon}(h|k) := 1 - \operatorname{err}_{\varepsilon}(h|k)$, where $\operatorname{err}_{\varepsilon}(h|k) := P_{X|k}(\exists x' \in \operatorname{Ball}(X; \varepsilon) \mid h(x') \neq k)$ is the adversarial robustness error of h on class k.

• Small $\operatorname{acc}_{\varepsilon}(h|k) \implies h$ is vulnerable to attacks on class k.

Distance to error set: $d(h|k) := \mathbb{E}_{P_{X|k}}[d(X, B(h, k))]$ denotes the average distance of a sample point of true label k, from the error set $B(h, k) := \{x \in \mathcal{X} \mid h(x) \neq k\}$ of samples assigned to another label.

• Small $d(h|k) \implies h$ is vulnerable to attacks on class k.

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Problem setup: notations

Standard accuracy: $\operatorname{acc}(h|k) := 1 - \operatorname{err}(h|k)$, where $\operatorname{err}(h|k) := P_{X|k}(h(X) \neq k)$ is the error of *h* on class *k*.

• Small $\operatorname{acc}(h|k) \implies h$ is inaccurate on class k.

■ Adversarial robustness accuracy: $\operatorname{acc}_{\varepsilon}(h|k) := 1 - \operatorname{err}_{\varepsilon}(h|k)$, where $\operatorname{err}_{\varepsilon}(h|k) := P_{X|k}(\exists x' \in \operatorname{Ball}(X; \varepsilon) \mid h(x') \neq k)$ is the adversarial robustness error of h on class k.

• Small $\operatorname{acc}_{\varepsilon}(h|k) \implies h$ is vulnerable to attacks on class k.

Distance to error set: $d(h|k) := \mathbb{E}_{P_{X|k}}[d(X, B(h, k))]$ denotes the average distance of a sample point of true label k, from the error set $B(h, k) := \{x \in \mathcal{X} \mid h(x) \neq k\}$ of samples assigned to another label.

• Small $d(h|k) \implies h$ is vulnerable to attacks on class k.

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

A motivating example (from [Tsipras '18])

Consider the following classification problem: Prediction target: $Y \sim \text{Bern}(1/2, \{\pm 1\})$ based on $p \ge 2$ explanatory variables $X := (X^1, X^2, \dots, X^p)$ given by

Robust feature: $X^1 \mid Y = +Y$ w.p 70% and -Y w.p. 30%.

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

A motivating example (from [Tsipras '18])

Consider the following classification problem: **Prediction target**: $Y \sim \text{Bern}(1/2, \{\pm 1\})$ based on $p \ge 2$

explanatory variables $X := (X^1, X^2, \dots, X^p)$ given by

Robust feature: $X^1 | Y = +Y$ w.p 70% and -Y w.p. 30%.

Non-robust features: $X^j | Y \sim \mathcal{N}(\eta Y, 1)$, for j = 2, ..., p, where $\eta \sim p^{-1/2}$ is a fixed scalar which controls the difficulty.

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

A motivating example (from [Tsipras '18])

Consider the following classification problem:

Prediction target: $Y \sim \text{Bern}(1/2, \{\pm 1\})$ based on $p \ge 2$ explanatory variables $X := (X^1, X^2, \dots, X^p)$ given by

Robust feature: $X^1 | Y = +Y$ w.p 70% and -Y w.p. 30%.

Non-robust features: $X^j | Y \sim \mathcal{N}(\eta Y, 1)$, for j = 2, ..., p, where $\eta \sim p^{-1/2}$ is a fixed scalar which controls the difficulty.

The linear classifier $h_{\text{lin}}(x) \equiv \text{sign}(w^T x)$ with $w = (0, 1/p, \dots, 1/p)$, where we allow ℓ_{∞} -perturbations of maximum size $\varepsilon \ge 2\eta$, solves the problem perfectly (100% accuracy) but its adversarial robustness is zero!

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

A motivating example (from [Tsipras '18])

Consider the following classification problem:

Prediction target: $Y \sim \text{Bern}(1/2, \{\pm 1\})$ based on $p \ge 2$ explanatory variables $X := (X^1, X^2, \dots, X^p)$ given by

Robust feature: $X^1 | Y = +Y$ w.p 70% and -Y w.p. 30%.

Non-robust features: $X^j | Y \sim \mathcal{N}(\eta Y, 1)$, for j = 2, ..., p, where $\eta \sim p^{-1/2}$ is a fixed scalar which controls the difficulty.

The linear classifier $h_{\text{lin}}(x) \equiv \text{sign}(w^T x)$ with $w = (0, 1/p, \dots, 1/p)$, where we allow ℓ_{∞} -perturbations of maximum size $\varepsilon \ge 2\eta$, solves the problem perfectly (100% accuracy) but its adversarial robustness is zero!

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Proof.

The standard accuracy of the classifier writes

$$\begin{aligned} \mathsf{acc}(h_{\mathsf{lin}}) &:= \mathbb{P}_{(X,Y)} \left(h_{\mathsf{lin}}(X) = Y \right) = \mathbb{P} \left(Y w^T X \ge 0 \right) \\ &= \mathbb{P}_Y \left(\left(Y/(p-1) \right) \sum_{j \ge 2} \mathcal{N}(\eta Y, 1) \ge 0 \right) \\ &= \mathbb{P} \left(\mathcal{N}(\eta, 1/(p-1)) \ge 0 \right) = \mathbb{P} \left(\mathcal{N}(0, 1/(p-1)) \ge -\eta \right) \\ &= \mathbb{P} \left(\mathcal{N}(0, 1/(p-1)) \le \eta \right) \ge 1 - e^{-(p-1)\eta^2/2}, \end{aligned}$$

which is $\geq 1 - \delta$ if $\eta \geq \sqrt{2\log(1/\delta)}/(p-1)$.

 $\blacksquare \implies h_{\text{lin}} \text{ is quasi-perfect!}$

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Proof.

The standard accuracy of the classifier writes

$$\begin{aligned} \operatorname{acc}(h_{\operatorname{lin}}) &:= \mathbb{P}_{(X,Y)} \left(h_{\operatorname{lin}}(X) = Y \right) = \mathbb{P} \left(Y w^T X \ge 0 \right) \\ &= \mathbb{P}_Y \left((Y/(p-1)) \sum_{j \ge 2} \mathcal{N}(\eta Y, 1) \ge 0 \right) \\ &= \mathbb{P} \left(\mathcal{N}(\eta, 1/(p-1)) \ge 0 \right) = \mathbb{P} \left(\mathcal{N}(0, 1/(p-1)) \ge -\eta \right) \\ &= \mathbb{P} \left(\mathcal{N}(0, 1/(p-1)) \le \eta \right) \ge 1 - e^{-(p-1)\eta^2/2}, \end{aligned}$$
which is $\ge 1 - \delta$ if $\eta \ge \sqrt{2 \log(1/\delta)/(p-1)}.$

$$\blacksquare \implies h_{\operatorname{lin}} \text{ is quasi-perfect!}$$

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Proof.

$$\begin{aligned} &\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) := \mathbb{P}_{(X,Y)} \left(Yh_{\operatorname{lin}}(X + \Delta x) \ge 0 \; \forall \|\Delta x\|_{\infty} \le \varepsilon \right) \\ &= \mathbb{P}_{(X,Y)} \left(\inf_{\|\Delta x\|_{\infty} \le \varepsilon} Yw^{T}(X + \Delta x) \ge 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \sup_{\|\Delta x\|_{\infty} \le \varepsilon} Yw^{T}\Delta x \ge 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \|Yw\|_{1} \ge 0 \right) = \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \ge 0 \right) \\ &= \mathbb{P}(\mathcal{N}(0, 1/(p-1)) \ge \varepsilon - \eta) \le e^{-(p-1)(\varepsilon - \eta)^{2}/2}. \end{aligned}$$
Thus $\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) \le \delta$ for $\varepsilon \ge \eta + \sqrt{2\log(1/\delta)/(p-1)}. \end{aligned}$

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Proof.

$$\begin{aligned} &\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) := \mathbb{P}_{(X,Y)} \left(Yh_{\operatorname{lin}}(X + \Delta x) \geq 0 \; \forall \|\Delta x\|_{\infty} \leq \varepsilon \right) \\ &= \mathbb{P}_{(X,Y)} \left(\inf_{\|\Delta x\|_{\infty} \leq \varepsilon} Yw^{T}(X + \Delta x) \geq 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \sup_{\|\Delta x\|_{\infty} \leq \varepsilon} Yw^{T}\Delta x \geq 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \|Yw\|_{1} \geq 0 \right) = \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \geq 0 \right) \\ &= \mathbb{P}(\mathcal{N}(0, 1/(p-1)) \geq \varepsilon - \eta) \leq e^{-(p-1)(\varepsilon - \eta)^{2}/2}. \end{aligned}$$
Thus $\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) \leq \delta$ for $\varepsilon \geq \eta + \sqrt{2\log(1/\delta)/(p-1)}. \end{aligned}$

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Proof.

$$\begin{aligned} &\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) := \mathbb{P}_{(X,Y)} \left(Yh_{\operatorname{lin}}(X + \Delta x) \geq 0 \; \forall \|\Delta x\|_{\infty} \leq \varepsilon \right) \\ &= \mathbb{P}_{(X,Y)} \left(\inf_{\|\Delta x\|_{\infty} \leq \varepsilon} Yw^{T}(X + \Delta x) \geq 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \sup_{\|\Delta x\|_{\infty} \leq \varepsilon} Yw^{T}\Delta x \geq 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \|Yw\|_{1} \geq 0 \right) = \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \geq 0 \right) \\ &= \mathbb{P}(\mathcal{N}(0, 1/(p-1)) \geq \varepsilon - \eta) \leq e^{-(p-1)(\varepsilon - \eta)^{2}/2}. \end{aligned}$$
Thus $\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) \leq \delta$ for $\varepsilon \geq \eta + \sqrt{2\log(1/\delta)/(p-1)}. \end{aligned}$

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Proof.

$$\begin{aligned} &\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) := \mathbb{P}_{(X,Y)} \left(Yh_{\operatorname{lin}}(X + \Delta x) \geq 0 \; \forall \|\Delta x\|_{\infty} \leq \varepsilon \right) \\ &= \mathbb{P}_{(X,Y)} \left(\inf_{\|\Delta x\|_{\infty} \leq \varepsilon} Yw^{T}(X + \Delta x) \geq 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \sup_{\|\Delta x\|_{\infty} \leq \varepsilon} Yw^{T}\Delta x \geq 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \|Yw\|_{1} \geq 0 \right) = \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \geq 0 \right) \\ &= \mathbb{P}(\mathcal{N}(0, 1/(p-1)) \geq \varepsilon - \eta) \leq e^{-(p-1)(\varepsilon - \eta)^{2}/2}. \end{aligned}$$
Thus $\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) \leq \delta$ for $\varepsilon \geq \eta + \sqrt{2\log(1/\delta)/(p-1)}. \end{aligned}$

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Proof.

the adversarial robustness accuracy of h_{lin} writes

$$\begin{aligned} &\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) := \mathbb{P}_{(X,Y)} \left(Yh_{\operatorname{lin}}(X + \Delta x) \ge 0 \; \forall \|\Delta x\|_{\infty} \le \varepsilon \right) \\ &= \mathbb{P}_{(X,Y)} \left(\inf_{\|\Delta x\|_{\infty} \le \varepsilon} Yw^{T}(X + \Delta x) \ge 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \sup_{\|\Delta x\|_{\infty} \le \varepsilon} Yw^{T}\Delta x \ge 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \|Yw\|_{1} \ge 0 \right) = \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \ge 0 \right) \\ &= \mathbb{P}(\mathcal{N}(0, 1/(p-1)) \ge \varepsilon - \eta) \le e^{-(p-1)(\varepsilon - \eta)^{2}/2}. \end{aligned}$$
Thus $\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) \le \delta$ for $\varepsilon \ge \eta + \sqrt{2\log(1/\delta)/(p-1)}. \end{aligned}$

That is, the **adversarial accuracy** of *h*_{lin} is close to **zero**!

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Proof.

the adversarial robustness accuracy of h_{lin} writes

$$\begin{aligned} &\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) := \mathbb{P}_{(X,Y)} \left(Yh_{\operatorname{lin}}(X + \Delta x) \ge 0 \; \forall \|\Delta x\|_{\infty} \le \varepsilon \right) \\ &= \mathbb{P}_{(X,Y)} \left(\inf_{\|\Delta x\|_{\infty} \le \varepsilon} Yw^{T}(X + \Delta x) \ge 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \sup_{\|\Delta x\|_{\infty} \le \varepsilon} Yw^{T}\Delta x \ge 0 \right) \\ &= \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \|Yw\|_{1} \ge 0 \right) = \mathbb{P}_{(X,Y)} \left(Yw^{T}X - \varepsilon \ge 0 \right) \\ &= \mathbb{P}(\mathcal{N}(0, 1/(p-1)) \ge \varepsilon - \eta) \le e^{-(p-1)(\varepsilon - \eta)^{2}/2}. \end{aligned}$$
Thus $\operatorname{acc}_{\varepsilon}(h_{\operatorname{lin}}) \le \delta$ for $\varepsilon \ge \eta + \sqrt{2\log(1/\delta)/(p-1)}. \end{aligned}$

That is, the adversarial accuracy of h_{lin} is close to zero!

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

What could be going on ? [Intuition from Tsipras and co.]

Prediction target: $Y \sim \text{Bern}(1/2, \{\pm 1\})$ Robust feature: $X^1 \mid Y = +Y$ w.p 70% and -Y w.p. 30%. Non-robust features: $X^j \mid Y \sim \mathcal{N}(\eta Y, 1)$, for j = 2, ..., p



BTW, we note that an optimal adversarial attack can be done by taking $\Delta x^1 = 0$ and $\Delta x^j = -\varepsilon y$ for all j = 2, ..., p.

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

What could be going on ? [Intuition from Tsipras and co.]

Prediction target: $Y \sim \text{Bern}(1/2, \{\pm 1\})$ Robust feature: $X^1 \mid Y = +Y$ w.p 70% and -Y w.p. 30%. Non-robust features: $X^j \mid Y \sim \mathcal{N}(\eta Y, 1)$, for j = 2, ..., p



BTW, we note that an optimal adversarial attack can be done by taking $\Delta x^1 = 0$ and $\Delta x^j = -\varepsilon y$ for all $j = 2, \dots, p$. Basic intuition: In standard training, all correlation is good correlation

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

What could be going on ? [Intuition from Tsipras and co.]

Prediction target: $Y \sim \text{Bern}(1/2, \{\pm 1\})$ Robust feature: $X^1 \mid Y = +Y$ w.p 70% and -Y w.p. 30%. Non-robust features: $X^j \mid Y \sim \mathcal{N}(\eta Y, 1)$, for j = 2, ..., p



BTW, we note that an optimal adversarial attack can be done by taking $\Delta x^1 = 0$ and $\Delta x^j = -\varepsilon y$ for all j = 2, ..., p. Basic intuition:

In standard training, all correlation is good correlation

If we want robustness, must avoid weakly correlated features

⇒ learn causal features ?

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

What could be going on ? [Intuition from Tsipras and co.]

Prediction target: $Y \sim \text{Bern}(1/2, \{\pm 1\})$ Robust feature: $X^1 \mid Y = +Y$ w.p 70% and -Y w.p. 30%. Non-robust features: $X^j \mid Y \sim \mathcal{N}(\eta Y, 1)$, for j = 2, ..., p



BTW, we note that an optimal adversarial attack can be done by taking $\Delta x^1 = 0$ and $\Delta x^j = -\varepsilon y$ for all j = 2, ..., p. Basic intuition:

In standard training, all correlation is good correlation

- If we want robustness, must avoid weakly correlated features
 - \implies learn causal features ?

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

BTW, humans are not "perfect"



Roberto Toro @R3RT0 · 23 avr. optical illusions are more like

Lionel Page @page_eco · 8 mars

Is it a duck or a rabbit?

Google Cloud Vision's algorithm has the same optical illusion than you and me. It sees one or the other, depending on how the image is rotated.



ht @minimaxir

Elvis Dohmatob

Limits on Robustness to Adversarial Examples - slide 15 / 41

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Talagrand transportation-cost inequality

The $T_2(c)$ property

Given $c \ge 0$, a distribution μ on \mathcal{X} is said to satisfy $T_2(c)$ if for every distribution ν on \mathcal{X} with $\nu \ll \mu$, one has

$$W_2(\nu,\mu) \le \sqrt{2c \operatorname{kl}(\nu \| \mu)},\tag{1}$$

where kl $(\nu \| \mu) := \int_{\mathcal{X}} \log(d\nu/d\mu) d\mu$, entropy of ν relative to μ .

Generalizes the well-known **Pinsker's inequality** for the total variation distance between probability measures (take 2c = 1/2).

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Talagrand transportation-cost inequality

The $T_2(c)$ property

Given $c \ge 0$, a distribution μ on \mathcal{X} is said to satisfy $T_2(c)$ if for every distribution ν on \mathcal{X} with $\nu \ll \mu$, one has

$$W_2(\nu,\mu) \le \sqrt{2c \operatorname{kl}(\nu \| \mu)},\tag{1}$$

where $kl(\nu \| \mu) := \int_{\mathcal{X}} log(d\nu/d\mu) d\mu$, entropy of ν relative to μ .

Generalizes the well-known **Pinsker's inequality** for the total variation distance between probability measures (take 2c = 1/2).

• Unlike Pinsker's inequality which holds unconditionally, the inequality $T_2(c)$ is a privilege only enjoyed by special classes of reference distributions μ .

Problem setup **No Free Lunch Theorems** The Strong No Free Lunch Theorem Corollaries

Talagrand transportation-cost inequality

The $T_2(c)$ property

Given $c \ge 0$, a distribution μ on \mathcal{X} is said to satisfy $T_2(c)$ if for every distribution ν on \mathcal{X} with $\nu \ll \mu$, one has

$$W_2(\nu,\mu) \le \sqrt{2c \operatorname{kl}(\nu \| \mu)},\tag{1}$$

where $kl(\nu \| \mu) := \int_{\mathcal{X}} log(d\nu/d\mu) d\mu$, entropy of ν relative to μ .

Generalizes the well-known **Pinsker's inequality** for the total variation distance between probability measures (take 2c = 1/2).

• Unlike Pinsker's inequality which holds unconditionally, the inequality $T_2(c)$ is a privilege only enjoyed by special classes of reference distributions μ .

Problem setup **No Free Lunch Theorems** The Strong No Free Lunch Theorem Corollaries

BLOWUP / aka concentration of measure

The BLOWUP(c) property

 μ is said to satisfy BLOWUP(c) if for every Borel $B \subseteq \mathcal{X}$ with $\mu(B) > 0$ and for every $\varepsilon \geq \sqrt{2c \log(1/\mu(B))}$, it holds that

$$\mu(B^{\varepsilon}) \ge 1 - e^{-\frac{1}{2c}(\varepsilon - \sqrt{2c\log(1/\mu(B))})^2}.$$
(2)

It is a classical result that the Gaussian distribution on \mathbb{R}^p has BLOWUP(1) and T₂(1), a phenomenon known as **Gaussian isoperimetry**.

Problem setup **No Free Lunch Theorems** The Strong No Free Lunch Theorem Corollaries

BLOWUP / aka concentration of measure

The BLOWUP(c) property

 μ is said to satisfy BLOWUP(c) if for every Borel $B \subseteq \mathcal{X}$ with $\mu(B) > 0$ and for every $\varepsilon \geq \sqrt{2c \log(1/\mu(B))}$, it holds that

$$\mu(B^{\varepsilon}) \ge 1 - e^{-\frac{1}{2c}(\varepsilon - \sqrt{2c\log(1/\mu(B))})^2}.$$
(2)

It is a classical result that the Gaussian distribution on \mathbb{R}^{p} has BLOWUP(1) and T₂(1), a phenomenon known as Gaussian isoperimetry.

 These result dates back to works of Borel, Lévy, Talagrand and of Marton (see [Boucheron '13] textbook)

Problem setup **No Free Lunch Theorems** The Strong No Free Lunch Theorem Corollaries

BLOWUP / aka concentration of measure

The BLOWUP(c) property

 μ is said to satisfy BLOWUP(c) if for every Borel $B \subseteq \mathcal{X}$ with $\mu(B) > 0$ and for every $\varepsilon \geq \sqrt{2c \log(1/\mu(B))}$, it holds that

$$\mu(B^{\varepsilon}) \ge 1 - e^{-\frac{1}{2c}(\varepsilon - \sqrt{2c\log(1/\mu(B))})^2}.$$
(2)

It is a classical result that the Gaussian distribution on \mathbb{R}^{p} has BLOWUP(1) and T₂(1), a phenomenon known as Gaussian isoperimetry.

• These result dates back to works of Borel, Lévy, Talagrand and of Marton (see [Boucheron '13] textbook)
Problem setup **No Free Lunch Theorems** The Strong No Free Lunch Theorem Corollaries

Marton's Blowup lemma

Lemma (Marton's BLOWUP Lemma)

On a metric space, it holds that $T_2(c) \subseteq BLOWUP(c)$.

Proof. Fact: $kl(\mu|_B||\mu) = log(1/\mu(B))$, where $\mu|_B(A) := \frac{\mu(A \cap B)}{\mu(B)}$

Problem setup **No Free Lunch Theorems** The Strong No Free Lunch Theorem Corollaries

Marton's Blowup lemma

Lemma (Marton's BLOWUP Lemma)

On a metric space, it holds that $T_2(c) \subseteq BLOWUP(c)$.

Proof. Fact: $kl(\mu|_B || \mu) = log(1/\mu(B))$, where $\mu|_B(A) := \frac{\mu(A \cap B)}{\mu(B)}$.

Thus $\varepsilon \leq W_2(\mu|_B, \mu_{\mathcal{X} \setminus B^{\varepsilon}}) \leq W_2(\mu|_B, \mu) + W_2(\mu|_{\mathcal{X} \setminus B^{\varepsilon}}, \mu)$

Problem setup **No Free Lunch Theorems** The Strong No Free Lunch Theorem Corollaries

Marton's Blowup lemma

Lemma (Marton's BLOWUP Lemma)

On a metric space, it holds that $T_2(c) \subseteq BLOWUP(c)$.

Proof. Fact: $kl(\mu|_B || \mu) = log(1/\mu(B))$, where $\mu|_B(A) := \frac{\mu(A \cap B)}{\mu(B)}$.

 $\mathsf{Thus}\; \varepsilon \leq \mathit{W}_2(\mu|_B, \mu_{\mathcal{X} \setminus B^\varepsilon}) \leq \mathit{W}_2(\mu|_B, \mu) + \mathit{W}_2(\mu|_{\mathcal{X} \setminus B^\varepsilon}, \mu)$

 $\leq \sqrt{2c \operatorname{kl}(\mu|_B \| \mu)} + \sqrt{2c \operatorname{kl}(\mu|_{\mathcal{X} \setminus B^{\varepsilon}} \| \mu)}$

Problem setup **No Free Lunch Theorems** The Strong No Free Lunch Theorem Corollaries

Marton's Blowup lemma

Lemma (Marton's BLOWUP Lemma)

On a metric space, it holds that $T_2(c) \subseteq BLOWUP(c)$.

Proof. Fact: $kl(\mu|_B || \mu) = log(1/\mu(B))$, where $\mu|_B(A) := \frac{\mu(A \cap B)}{\mu(B)}$.

Thus
$$\varepsilon \leq W_2(\mu|_B, \mu_{\mathcal{X} \setminus B^\varepsilon}) \leq W_2(\mu|_B, \mu) + W_2(\mu|_{\mathcal{X} \setminus B^\varepsilon}, \mu)$$

$$\leq \sqrt{2c\,\mathsf{kl}(\mu|_{B}\|\mu)} + \sqrt{2c\,\mathsf{kl}(\mu|_{\mathcal{X}\setminus B^arepsilon}\|\mu)}$$

 $\leq \sqrt{2c\log(1/\mu(B))} + \sqrt{2c\log(1/\mu(\mathcal{X}\setminus B^{\varepsilon}))}$

Problem setup **No Free Lunch Theorems** The Strong No Free Lunch Theorem Corollaries

Marton's Blowup lemma

Lemma (Marton's BLOWUP Lemma)

On a metric space, it holds that $T_2(c) \subseteq BLOWUP(c)$.

Proof. Fact: $kl(\mu|_B || \mu) = log(1/\mu(B))$, where $\mu|_B(A) := \frac{\mu(A \cap B)}{\mu(B)}$.

Thus
$$\varepsilon \leq W_2(\mu|_B, \mu_{\mathcal{X} \setminus B^{\varepsilon}}) \leq W_2(\mu|_B, \mu) + W_2(\mu|_{\mathcal{X} \setminus B^{\varepsilon}}, \mu)$$

$$\leq \sqrt{2c \operatorname{kl}(\mu|B\|\mu)} + \sqrt{2c \operatorname{kl}(\mu|X \setminus B^{\varepsilon}\|\mu)}$$

$$\leq \sqrt{2c \log(1/\mu(B))} + \sqrt{2c \log(1/\mu(X \setminus B^{\varepsilon}))}$$

$$= \sqrt{2c \log(1/\mu(B))} + \sqrt{2c \log(1/(1-\mu(B^{\varepsilon})))}$$

Problem setup **No Free Lunch Theorems** The Strong No Free Lunch Theorem Corollaries

Marton's Blowup lemma

Lemma (Marton's BLOWUP Lemma)

On a metric space, it holds that $T_2(c) \subseteq BLOWUP(c)$.

Proof. Fact:
$$\mathsf{kl}(\mu|_B \| \mu) = \mathsf{log}(1/\mu(B))$$
 , where $\mu|_B(A) := rac{\mu(A \cap B)}{\mu(B)}$.

Thus
$$\varepsilon \leq W_2(\mu|_B, \mu_{\mathcal{X} \setminus B^{\varepsilon}}) \leq W_2(\mu|_B, \mu) + W_2(\mu|_{\mathcal{X} \setminus B^{\varepsilon}}, \mu)$$

$$egin{aligned} &\leq \sqrt{2c\,\mathsf{kl}(\mu|_{B}\|\mu)} + \sqrt{2c\,\mathsf{kl}(\mu|_{\mathcal{X}\setminus B^arepsilon}\|\mu)} \ &\leq \sqrt{2c\log(1/\mu(B))} + \sqrt{2c\log(1/\mu(\mathcal{X}\setminus B^arepsilon))} \ &= \sqrt{2c\log(1/\mu(B))} + \sqrt{2c\log(1/(1-\mu(B^arepsilon)))}. \end{aligned}$$

Rearranging the above inequality gives

$$\sqrt{2c\log(1/(1-\mu(B^{\varepsilon})))} \ge (\varepsilon - \sqrt{2c\log(1/\mu(B))})_+,$$

and the result follows after squaring & exponentiating.

Problem setup **No Free Lunch Theorems** The Strong No Free Lunch Theorem Corollaries

Marton's Blowup lemma

Lemma (Marton's BLOWUP Lemma)

On a metric space, it holds that $T_2(c) \subseteq BLOWUP(c)$.

Proof. Fact:
$$\mathsf{kl}(\mu|_B \| \mu) = \mathsf{log}(1/\mu(B))$$
 , where $\mu|_B(A) := rac{\mu(A \cap B)}{\mu(B)}$.

Thus
$$\varepsilon \leq W_2(\mu|_B, \mu_{\mathcal{X} \setminus B^{\varepsilon}}) \leq W_2(\mu|_B, \mu) + W_2(\mu|_{\mathcal{X} \setminus B^{\varepsilon}}, \mu)$$

$$\begin{split} &\leq \sqrt{2c\,\mathsf{kl}(\mu|_B\|\mu)} + \sqrt{2c\,\mathsf{kl}(\mu|_{\mathcal{X}\setminus B^\varepsilon}\|\mu)} \\ &\leq \sqrt{2c\,\mathsf{log}(1/\mu(B))} + \sqrt{2c\,\mathsf{log}(1/\mu(\mathcal{X}\setminus B^\varepsilon))} \\ &= \sqrt{2c\,\mathsf{log}(1/\mu(B))} + \sqrt{2c\,\mathsf{log}(1/(1-\mu(B^\varepsilon)))}. \end{split}$$

Rearranging the above inequality gives

$$\sqrt{2c\log(1/(1-\mu(B^{\varepsilon})))} \geq (\varepsilon - \sqrt{2c\log(1/\mu(B))})_+,$$

and the result follows after squaring & exponentiating.

Problem setup No Free Lunch Theorems **The Strong No Free Lunch Theorem** Corollaries

Adversarial attacks are a 'butterfly effect' on data manifold

Error set: $B(h, k) = \{x \in \mathcal{X} \mid h(x) \neq k\}, h = \text{classifier}$

• Neighbors of error set: $B(h,k)^{\varepsilon} := \{x \in \mathcal{X} \mid d(x,B(h,k)) \leq \varepsilon\}$



err(h|k) := $P_{X|k}(B(h,k)) > 0$ if h is not perfect on class k.
Consequence is that $\operatorname{acc}_{\varepsilon}(h|k) \searrow 0$ expo. fast as function of ε .
Thus adversarial robustness is impossible in general!
Manuscript: https://arxiv.org/pdf/1810.04065.pdf

Problem setup No Free Lunch Theorems **The Strong No Free Lunch Theorem** Corollaries

Strong No Free Lunch Theorem

Theorem (Strong "No Free Lunch" [Dohmatob '18])

Suppose that conditional distribution $P_{X|k}$ has the $T_2(\sigma_k^2)$ property. Given a classifier $h : \mathcal{X} \mapsto \mathcal{Y}$ such that $\operatorname{err}(h|k) > 0$, define $\varepsilon(h|k) := \sigma_k \sqrt{2\log(1/\operatorname{err}(h|k))}$. Then we have the following bounds:

(A) Adversarial robustness accuracy: if $\varepsilon \ge \varepsilon(h|k)$, then

$$\operatorname{acc}_{\varepsilon}(h|k) \leq e^{-rac{1}{2\sigma_k^2}(\varepsilon - \varepsilon(h|k))^2}.$$
 (3)

(B) Average distance to error set:

$$d(h|k) \le \sigma_k \left(\sqrt{\log(1/\operatorname{err}(h|k))} + \sqrt{\pi/2} \right)$$
(4)

Problem setup No Free Lunch Theorems **The Strong No Free Lunch Theorem** Corollaries

Proof



• Use Marton's Lemma: BLOWUP $(\sigma_k^2) \subseteq \mathsf{T}_2(\sigma_k^2)$ with $B := B(h, k) := \{x \in \mathcal{X} \mid h(x) \neq k\}$ and $\mu = P_{X|k}$.

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem **Corollaries**

Corollary (Strong "No Free Lunch" Theorem on flat space)

Let $1 \le q \le \infty$. Define $\epsilon_q(h|k) := \varepsilon(h|k)p^{1/q-1/2}$. If in addition to the assumptions of Strong No Free Lunch Theorem, and suppose the feature space is flat, i.e $\operatorname{Ric}_{\mathcal{X}} = 0$, then for the ℓ_q threat model, we have the following bounds:

(A1) Adversarial robustness accuracy: if $\varepsilon \ge \epsilon_q(h|k)$, then

$$\operatorname{acc}_{\varepsilon}(h|k) \leq e^{-rac{p^{1-2/q}}{2\sigma_{k}^{2}}(\varepsilon-\epsilon_{q}(h|k))^{2}}.$$
 (5)

(A2) Average distance to error set:

$$d(h|k) \leq \sigma_k p^{1/q-1/2} \left(\sqrt{\log(1/\operatorname{err}(h|k))} + \sqrt{\pi/2} \right).$$
 (6)

Note that the case q = 1 is a proxy for "few-pixel" attack models [Su et a. '18].

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem **Corollaries**

Strong No Free Lunch Theorem

Corollary (Strong NFLT for ℓ_{∞} attacks [Dohmatob '18])

In particular, for the ℓ_∞ threat model, we have the following bounds:

(B1) Adversarial robustness accuracy: If $\varepsilon \ge \varepsilon(h|k)/\sqrt{p}$, then

$$\operatorname{acc}_{\varepsilon}(h|k) \leq e^{-rac{p}{2\sigma_{k}^{2}}(\varepsilon - \varepsilon(h|k)/\sqrt{p})^{2}}.$$
 (7)

(B2) Average distance to error set:

$$d(h|k) \leq \frac{\sigma_k}{\sqrt{p}} \left(\sqrt{\log(1/\operatorname{err}(h|k))} + \sqrt{\pi/2} \right)$$
(8)

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem **Corollaries**

Special cases of our results

■ Log-concave distribs $dP_{X|k} \propto e^{-\nu_k(x)} dx$ satisfying Emery-Bakry curvature condition: $\operatorname{Hess}_x(\nu_k) + \operatorname{Ric}_x(\mathcal{X}) \succeq (1/\sigma_k^2) I_p$.

• e.g multi-variate Gaussian (considered in [Tsipras '18, Fawzi et al. 18])

Perturbed log-concave distribs (via Holley-Shroock Theorem)

The uniform measure on compact Riemannian manifolds of positive Ricci curvature, e.g spheres (considered in [Gilmer '18]), tori, or any compact Lie group.

Pushforward via a Lipschitz function f, of a distribution in $T_2(\sigma_k^2)$. Indeed, take $\tilde{\sigma}_k = \|f\|_{\text{Lip}}\sigma_k$.

etc.

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem **Corollaries**

Worked example: Adversarial spheres [Gilmer '18]



• $Y \sim \text{Bern}(1/2, \{\pm\})$,

 $X|k \sim uniform(\mathbb{S}_{R_k}^p)$, where $R_+ > R_- > 0$.

■ $\mathbb{S}_{R_k}^p$ is a compact Riemannian manifold with constant Ricci curvature $(p-1)R_k^{-2}$.

Thus $P_{X|k}$ satisfies $T_2(R_k^2/(p-1))$.

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem **Corollaries**

Worked example: Adversarial spheres [Gilmer '18]



- $Y \sim \text{Bern}(1/2, \{\pm\}),$
- $X|k \sim uniform(\mathbb{S}_{R_k}^p)$, where $R_+ > R_- > 0$.
- $\mathbb{S}_{R_k}^{p}$ is a compact Riemannian manifold with constant Ricci curvature $(p-1)R_k^{-2}$.

Thus
$$P_{X|k}$$
 satisfies $T_2(R_k^2/(p-1))$.

$$\therefore \mathbb{E}_{X|k}[d_{\text{geo}}(X, B(h, k))] \leq \frac{R_k}{\sqrt{p-1}} (\sqrt{2\log(1/\operatorname{err}(h|k))} + \sqrt{\pi/2}) \\ \sim \frac{R_k}{\sqrt{p}} \Phi^{-1}(\operatorname{acc}(h|k))$$

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Worked example: Adversarial spheres [Gilmer '18]



- $Y \sim \text{Bern}(1/2, \{\pm\})$,
- $X|k \sim uniform(\mathbb{S}^p_{R_k})$, where $R_+ > R_- > 0$.

■ $\mathbb{S}_{R_k}^{p}$ is a compact Riemannian manifold with constant Ricci curvature $(p-1)R_k^{-2}$.

Thus
$$P_{X|k}$$
 satisfies $T_2(R_k^2/(p-1))$.

$$egin{aligned} & \therefore \mathbb{E}_{X|k}[d_{ ext{geo}}(X,B(h,k))] \leq rac{R_k}{\sqrt{p-1}}(\sqrt{2\log(1/err(h|k))} + \sqrt{\pi/2}) \ & \sim rac{R_k}{\sqrt{p}} \Phi^{-1}(\operatorname{acc}(h|k)) \end{aligned}$$

This is the same bound obtained in "manually" [Gilmer '18].

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem **Corollaries**

Worked example: Adversarial spheres [Gilmer '18]



- $Y \sim \text{Bern}(1/2, \{\pm\}),$
- $X|k \sim uniform(\mathbb{S}^p_{R_k})$, where $R_+ > R_- > 0$.

■ $\mathbb{S}_{R_k}^{p}$ is a compact Riemannian manifold with constant Ricci curvature $(p-1)R_k^{-2}$.

Thus
$$P_{X|k}$$
 satisfies $T_2(R_k^2/(p-1))$.

$$egin{aligned} & \therefore \mathbb{E}_{X|k}[d_{ ext{geo}}(X,B(h,k))] \leq rac{R_k}{\sqrt{p-1}}(\sqrt{2\log(1/err(h|k))} + \sqrt{\pi/2}) \ & \sim rac{R_k}{\sqrt{p}} \Phi^{-1}(\operatorname{acc}(h|k)) \end{aligned}$$

This is the same bound obtained in "manually" [Gilmer '18].

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem **Corollaries**

Some empirical confirmation



Phase-transition occurs as predicted by our theorems

Problem setup No Free Lunch Theorems The Strong No Free Lunch Theorem Corollaries

Key papers

- [Tsipras '18] There is no free lunch in adversarial robustness
- ■[Gilmer '18] Adversarial spheres
- [Fawzi '18] Adversarial vulnerability for any classifier
- [Athalye '18] Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples
- [Dohmatob '19] Generalized No Free Lunch Theorem for Adversarial Robustness
- [Shafahi '19] Are adversarial examples inevitable?

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

Universal lower bounds

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

Overview

Previous classifier-dependent bounds make very strong assumptions on both the data and the classifier (e.g the theory fails for perfect classifiers)

It would be nice to have **universal** bounds which only depend on the geomtry of the class-conditional distributions P_+ and P_-

This is very very recent work, started by [Bhagoji '19] (to appear in NeurIPS!)

My own work builds on [Bhagoji '19] as is still largely ongoing (AISTATS ???)

References

 [Bhagoji '19] Lower Bounds on Adversarial Robustness from Optimal Transport

Abstract view of adversarial attacks

- The feature space \mathcal{X} is an abstract measure space, and the target space is $\{\pm 1\}$ (binary classification). E.g $\mathcal{X} = (\mathbb{R}^p, \text{Borell})$.
- •Let *P* be an unknown probability distribution on the product space $\mathcal{X} \times \{\pm 1\}$.
- A classifier is any measurable function $h : \mathcal{X} \to \{\pm 1\}$.
- •An attack-model \mathcal{A} is the prescription of a closed neighborhood \mathcal{A}_x for each point x of \mathcal{X} . E.g $\mathcal{A}_x = \text{Ball}_{\ell_{\infty}}(x;\varepsilon)$. The case $\mathcal{A}_x = \{x\} \ \forall x \in \mathcal{X}$ corresponds to the attackless model.

■A type- \mathcal{A} attack is any measurable function $a : \mathcal{X} \times \{\pm 1\} \to \mathcal{X}$ such that $a(x, y) \in \mathcal{A}_x \forall (x, y) \in \mathcal{X} \times \{\pm 1\}$. With abuse of notation, we'll also write $a \in \mathcal{A}$. E.g a(x, y) := x - yz for some fixed $z \in \text{Ball}_{\ell_{\infty}}(0; \varepsilon)$.

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

Abstract view of adversarial attacks

The robustness error of h to type-A attacks is

$$\operatorname{err}_{\mathcal{A}}(h) := \mathbb{E}_{(x,y)\sim P}[\exists x' \in \mathcal{A}_x \text{ s.t } h(x') \neq y]$$
 (9)

The Bayes-optimal robustness error for type-A attacks is

$$\operatorname{err}_{\mathcal{A}}^* := \inf_{h} \operatorname{err}_{\mathcal{A}}(h)$$
 (10)

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

The flying pig example again!



(Picture is courtesy of https://gradientscience.org/intro_adversarial/)

■ $x \mapsto x^{\text{adv}} := x + \text{ noise}$, $\|\text{noise}\| \le \varepsilon = 0.005$ (in example above) ■ $\mathcal{X} = \mathbb{R}^{\#\text{pixels}}$, $\mathcal{A}_x = \text{Ball}_{\ell_{\infty}}(x; 0.005)$

Adversarial attacks as optimal transport [Bhagoji '19]

Given a classifier h, consider the derived classifier $\tilde{h}: \mathcal{X} \to \{\pm 1\}$

$$\tilde{h}(x) := \begin{cases} y, & \text{if } \exists y \in \{\pm 1\} \text{ s.t } h(x') = y \ \forall x' \in \mathcal{A}_x, \\ \bot, & \text{else.} \end{cases}$$
(11)

Define the transport ground-cost



$$c_{\mathcal{A}}(x,x') = egin{cases} 1, & ext{if } \mathcal{A}_x \cap \mathcal{A}_{x'} = \emptyset, \ 0, & ext{else}, \end{cases}$$

Adversarial attacks as optimal transport [Bhagoji '19]

Given a classifier h, consider the derived classifier $\tilde{h}: \mathcal{X} \to \{\pm 1\}$

$$\tilde{h}(x) := \begin{cases} y, & \text{if } \exists y \in \{\pm 1\} \text{ s.t } h(x') = y \ \forall x' \in \mathcal{A}_x, \\ \bot, & \text{else.} \end{cases}$$
(11)



Define the transport ground-cost

$$c_{\mathcal{A}}(x,x') = egin{cases} 1, & ext{if } \mathcal{A}_x \cap \mathcal{A}_{x'} = \emptyset, \ 0, & ext{else}, \end{cases}$$

and note that $\forall x, x' \in \mathcal{X}$, one has $\mathbb{1}_{\{\tilde{h}(x)=1\}} + \mathbb{1}_{\{\tilde{h}(x)=-1\}} \leq c_{\mathcal{A}}(x, x') + 1$,

i.e
$$f(x) - g(x') \leq c_{\mathcal{A}}(x, x')$$
.

Adversarial attacks as optimal transport [Bhagoji '19]

Given a classifier h, consider the derived classifier $\tilde{h}: \mathcal{X} \to \{\pm 1\}$

П

$$\tilde{h}(x) := \begin{cases} y, & \text{if } \exists y \in \{\pm 1\} \text{ s.t } h(x') = y \ \forall x' \in \mathcal{A}_x, \\ \bot, & \text{else.} \end{cases}$$
(11)



Define the transport ground-cost

$$c_{\mathcal{A}}(x,x') = egin{cases} 1, & ext{if } \mathcal{A}_x \cap \mathcal{A}_{x'} = \emptyset, \ 0, & ext{else}, \end{cases}$$

and note that $\forall x, x' \in \mathcal{X}$, one has $\mathbb{1}_{\{\tilde{h}(x)=1\}} + \mathbb{1}_{\{\tilde{h}(x)=-1\}} \leq c_{\mathcal{A}}(x, x') + 1,$

i.e
$$f(x) - g(x') \leq c_{\mathcal{A}}(x, x')$$
.

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

Adversarial attacks as optimal transport [Bhagoji '19]



and so (f_h, g_h) is a pair of Kantorovich potentials for OT with ground-cost c_A .

$$\therefore OT_{c_{\mathcal{A}}}(P_{-}, P_{+}) := \sup_{K-\text{potentials } \phi, \psi} \mathbb{E}_{P_{-}}[\phi(x)] - \mathbb{E}_{P_{+}}[\psi(x)]$$

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

Adversarial attacks as optimal transport [Bhagoji '19]



$$lacksymbol{\sigma} c_{\mathcal{A}}(x,x') = egin{cases} 1, & ext{if } \mathcal{A}_x \cap \mathcal{A}_{x'} = \emptyset, \ 0, & ext{else}, \end{cases}$$

■ $f_h(x) - g_h(x') \le c_A(x, x') \forall x, x' \in \mathcal{X}$, and so (f_h, g_h) is a pair of Kantorovich potentials for OT with ground-cost c_A .

$$\therefore OT_{c_{\mathcal{A}}}(P_{-}, P_{+}) := \sup_{\substack{K - \text{potentials } \phi, \psi \\ b \in \mathbb{E}_{P_{-}}[g_{h}(x)] - \mathbb{E}_{P_{+}}[f_{h}(x)]} \mathbb{E}_{P_{-}}[g_{h}(x)] - \mathbb{E}_{P_{+}}[f_{h}(x)]$$

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

Adversarial attacks as optimal transport [Bhagoji '19]



$$lacksymbol{\sigma} c_{\mathcal{A}}(x,x') = egin{cases} 1, & ext{if } \mathcal{A}_x \cap \mathcal{A}_{x'} = \emptyset, \ 0, & ext{else}, \end{cases}$$

■ $f_h(x) - g_h(x') \le c_A(x, x') \forall x, x' \in \mathcal{X}$, and so (f_h, g_h) is a pair of Kantorovich potentials for OT with ground-cost c_A .

$$\therefore OT_{c_{\mathcal{A}}}(P_{-}, P_{+}) := \sup_{\substack{K - \text{potentials } \phi, \psi \\ h \in \mathbb{F}_{P_{-}}[g_{h}(x)] - \mathbb{E}_{P_{+}}[f_{h}(x)]} \\ \ge \sup_{\substack{h \\ h}} \mathbb{E}_{P_{-}}[g_{h}(x)] - \mathbb{E}_{P_{+}}[f_{h}(x)] \\ = \sup_{\substack{h \\ h}} 2(1 - \text{err}_{\mathcal{A}}(h)) - 1 = 1 - \text{err}_{\mathcal{A}}^{*}$$

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

Adversarial attacks as optimal transport [Bhagoji '19]



$$lacksymbol{\sigma} c_{\mathcal{A}}(x,x') = egin{cases} 1, & ext{if } \mathcal{A}_x \cap \mathcal{A}_{x'} = \emptyset, \ 0, & ext{else}, \end{cases}$$

■ $f_h(x) - g_h(x') \le c_A(x, x') \forall x, x' \in \mathcal{X}$, and so (f_h, g_h) is a pair of Kantorovich potentials for OT with ground-cost c_A .

$$\therefore OT_{c_{\mathcal{A}}}(P_{-}, P_{+}) := \sup_{\substack{K - \text{potentials } \phi, \psi \\ h \in \mathbb{P}_{-}[g_{h}(x)] - \mathbb{E}_{P_{+}}[f_{h}(x)]} \\ \ge \sup_{h} \mathbb{E}_{P_{-}}[g_{h}(x)] - \mathbb{E}_{P_{+}}[f_{h}(x)] \\ = \sup_{h} 2(1 - \text{err}_{\mathcal{A}}(h)) - 1 = 1 - \text{err}_{\mathcal{A}}^{*}$$

Universal lower bound on adversarial robustness error

Theorem ([Bhagoji '19])

Given an attack model A, let $OT_A(P_+, P_-)$ be the **optimal transport** distance between the +ve and -ve class-conditonal distributions of the samples, with the ground cost given by $c_A(x, x') = \mathbb{1}_{\{A_x \cap A_{x'} = \emptyset\}}$. Then we have he following lower bound on the classification error against A-attacks

$$\operatorname{err}_{\mathcal{A}}^* \geq \frac{1}{2}(1 - OT_{\mathcal{A}}(P_+, P_-))$$
 (12)

In particular, for the attackless case where $A_x = \{x\} \forall x \in \mathcal{X}$, one has $c_{\mathcal{A}}(x, x') = \mathbb{1}_{x \neq x'}$ and so $OT_{\mathcal{A}}(P_+, P_-) = TV(P_+, P_-)$. The theorem then reduces to the following well-known result

$$\operatorname{err}^* \geq \frac{1}{2}(1 - TV(P_+, P_-)).$$
 (13)

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

Total-Variational reformulation of the bounds

Theorem ([Dohmatob '20 ?] / ongoing work)

• Let \mathcal{A} be an attack, and for $a \in \mathcal{A}$ define $a_+(x) :\equiv a(x,+1)$. Define $\Omega := \{(x,x') \in \mathcal{X}^2 \mid \mathcal{A}_x \cap \mathcal{A}_{x'} \neq \emptyset\}$, and

$$TV_{\mathcal{A}}(P_{-}, P_{+}) := \inf_{a \in \mathcal{A}} TV(a_{-\#}P_{-}, a_{+\#}P_{+}),$$

$$\widetilde{TV}_{\mathcal{A}}(P_{-}, P_{+}) := \inf_{\gamma_{1}, \gamma_{2}} TV(\operatorname{proj}_{2\#}\gamma_{1}, \operatorname{proj}_{1\#}\gamma_{2}),$$
(14)

where the inf is taken over all distributions on \mathcal{X}^2 which are concentrated on Ω s.t proj_{1#} $\gamma_1 = P_-$ and proj_{2#} $\gamma_2 = P_+$. Then,

$$OT_{\mathcal{A}}(P_{-},P_{+}) = \widetilde{TV}_{\mathcal{A}}(P_{-},P_{+}) \le TV_{\mathcal{A}}(P_{-},P_{+}),$$
(15)

and there is equality if P_{-} and P_{+} have densities w.r.t Lebesgue.

Above bound suggest that rather than doing adversarial training, we'd rather do normal training on adversarially augmented data!

Worked example: hierarchical Gaussian classification

(Example from [Schmidt '18]) $\mu \sim \mathcal{N}(0, I_p), y \sim \text{Bern}(\{\pm 1\}), X|(Y = y) \sim \mathcal{N}(y\mu, \sigma^2 I_p),$

• Consider the ℓ_{∞} -norm attack model \mathcal{A} give by $\mathcal{A}_{x} = \mathsf{Ball}_{\ell_{\infty}}(x;\varepsilon)$.

Given *n* samples $S_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ from this model, how small can the robust error of a classifier be ?

More precisely, lets bound

$$\mathbb{E}_{\mu \sim \mathcal{N}(0,1)} \inf_{\hat{h}} \mathbb{E}_{S_n \sim P^n} \mathbb{E}_{\hat{h}_n \sim \hat{h}(S_n)}[\operatorname{err}_{\mathcal{A}}(\hat{h}_n; \mu)],$$
(16)

where $\operatorname{err}_{\mathcal{A}}(\hat{h}_n; \mu)$ is the adversarial robust error of \hat{h}_n defined by $\operatorname{err}_{\mathcal{A}}(\hat{h}_n; \mu) := \mathbb{E}_{y \sim \operatorname{Bern}(\{\pm 1\})} \mathbb{E}_{x \sim \mathcal{N}(y\mu,\sigma^2 I_p)}[\exists x' \in \mathcal{A}_x \text{ s.t } \hat{h}_n(x') \neq y].$

Worked example: hierarchical Gaussian classification

$$\mathsf{err}_{\mathcal{A}}(\hat{h}_n;\mu) := \mathbb{E}_{y \sim \mathsf{Bern}(\{\pm 1\})} \mathbb{E}_{x \sim \mathcal{N}(y\mu,\sigma^2 I_p)} [\exists x' \in \mathcal{A}_x \text{ s.t } \hat{h}_n(x') \neq y].$$

The posterior distribution of the model parameter is $\mathcal{N}(\hat{\mu}_n, \hat{\sigma}_n^2)$, with $\hat{\sigma}_n^2 = \frac{\sigma^2}{\sigma^2 + n}$, and $\hat{\mu}_n = \frac{n}{\sigma^2 + n} \bar{x}$ with $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ $\therefore \operatorname{err}_{\mathcal{A}}(\hat{h}_n) = \inf_{h} \operatorname{err}_{\mathcal{A}}(h; \hat{\mu}_n, \hat{\sigma}_n^2) \geq \dots$ $\geq \mathbb{E}_{\mu} \mathbb{E}_{S_n} \frac{1}{2} \left(1 - \inf_{\|z\|_{\infty} \leq \varepsilon} \mathsf{TV}(\mathcal{N}(-\hat{\mu}_n + z, \hat{\sigma}_n^2), \mathcal{N}(\hat{\mu}_n - z, \hat{\sigma}_n^2)) \right)$ $\geq \mathbb{E}_{\mu}\mathbb{E}_{S_n}\Phi\left(\frac{\sqrt{p}}{\hat{\sigma}_{\pi}}(\|\hat{\mu}_n\|-\varepsilon)_+\right) \geq \mathbb{E}_{\mu}\mathbb{E}_{S_n}\mathbb{P}(\|\hat{\mu}_n\|_{\infty} \leq \varepsilon)\Phi(0)$ $\approx \frac{1}{2} \mathbb{P}_{u \sim \mathcal{N}(0, l_p)} \left(\frac{n}{n + \sigma^2} \| u \|_{\infty} \leq \varepsilon \right)$

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

Worked example: hierarchical Gaussian classification

Thus, if
$$n \leq \frac{\varepsilon^2 \sigma^2}{8 \log(d)}$$
, then $\frac{n}{\sigma^2 + n} \leq \frac{\varepsilon}{2\sqrt{2 \log(d)}}$, and so
 $\therefore \operatorname{err}_{\mathcal{A}, n}^* \geq \frac{1}{2} \mathbb{P}_{u \sim \mathcal{N}(0, l_p)} \left(\frac{n}{n + \sigma^2} \|u\|_{\infty} \leq \varepsilon \right)$
 $\geq \frac{1}{2} \mathbb{P}_{u \sim \mathcal{N}(0, l_p)} \left(\|u\|_{\infty} \leq 2\sqrt{2 \log(d)} \right)$
 \vdots
 $\geq \frac{1}{2} (1 - 1/d) \approx \frac{1}{2}$ in high dimensions
Preliminaries on adversarial robustness Classifier-dependent lower bounds Universal lower bounds

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

Main references

[Dohmatob '19] Generalized No Free Lunch Theorem for Adversarial Robustness

- **[**Bhagoji '19] Lower Bounds on Adversarial Robustness from Optimal Transport
- [Tsipras '18] There is no free lunch in adversarial robustness
 [Gilmer '18] Adversarial spheres
- **[**Goodfellow '14] *Explaining and harnessing adversarial examples*
- [Su '17] One pixel attack for fooling deep neural networks
- [Fawzi '18] Adversarial vulnerability for any classifier
- [Athalye '18] Obfuscated Gradients Give a False Sense of
- Security: Circumventing Defenses to Adversarial Examples

Preliminaries on adversarial robustness Classifier-dependent lower bounds Universal lower bounds

Link between adversarial examples and optimal transport Adversarially robust learning via adversarially augmented data

Questions ?

