

# Regularization properties of adversarially-trained linear regression

---

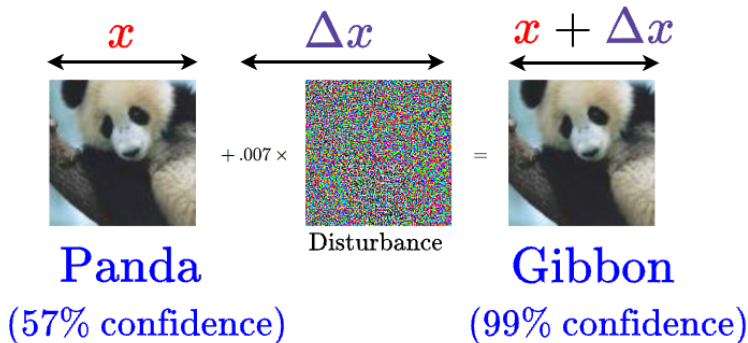
Antônio H. Ribeiro<sup>1,\*</sup>, Dave Zachariah<sup>1</sup>,  
Francis Bach<sup>2</sup>, Thomas B. Schon<sup>1</sup>

<sup>1</sup>Uppsala University, Sweden

<sup>2</sup>INRIA / PSL research university, France

\*Presenting

# Adversarial attacks

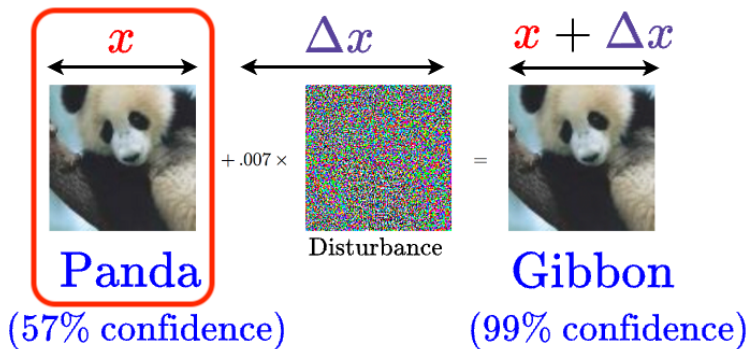


Explaining and Harnessing Adversarial Examples

I. J. Goodfellow, J. Shlens, C. Szegedy

ICLR (2015)

# Adversarial attacks

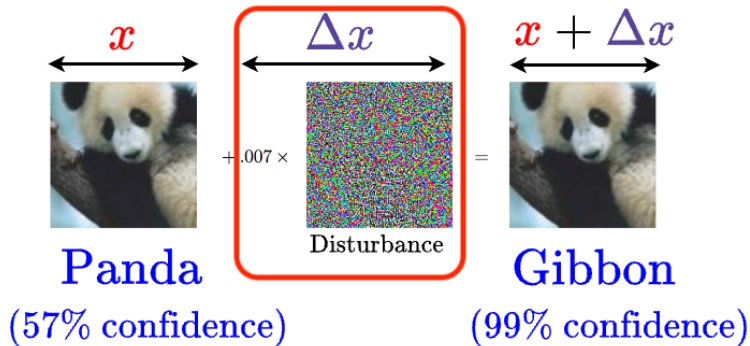


Explaining and Harnessing Adversarial Examples

I. J. Goodfellow, J. Shlens, C. Szegedy

ICLR (2015)

# Adversarial attacks

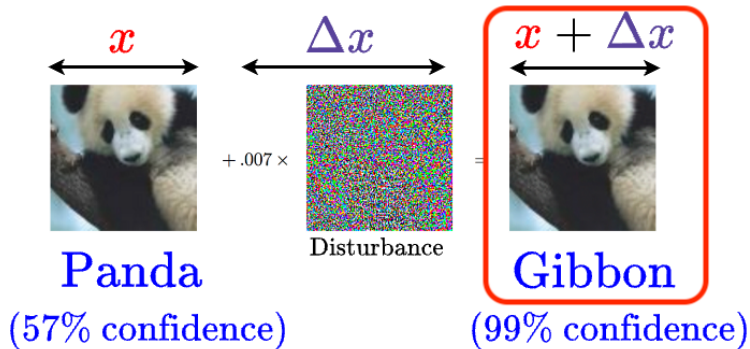


Explaining and Harnessing Adversarial Examples

I. J. Goodfellow, J. Shlens, C. Szegedy

ICLR (2015)

# Adversarial attacks



Explaining and Harnessing Adversarial Examples

I. J. Goodfellow, J. Shlens, C. Szegedy

ICLR (2015)

**Adversarial training:** *Each training sample is modified by an adversary.*

# Adversarially-trained linear regression

► **Linear regression:**

$$\min_{\beta} \sum_{i=1}^{\#train} (y_i - \beta^T x_i)^2$$

# Adversarially-trained linear regression

► **Linear regression:**

$$\min_{\beta} \sum_{i=1}^{\#train} \left( \underbrace{y_i}_{\text{observed}} - \underbrace{\beta^T x_i}_{\text{linear prediction}} \right)^2$$



# Adversarially-trained linear regression

- ▶ **Linear regression:**

$$\min_{\beta} \sum_{i=1}^{\#train} (y_i - \beta^T x_i)^2$$

- ▶ **Adversarial training** in linear regression:

$$(y_i - \beta^T (x_i + \Delta x_i))^2$$

# Adversarially-trained linear regression

- ▶ **Linear regression:**

$$\min_{\beta} \sum_{i=1}^{\#train} (y_i - \beta^T x_i)^2$$

- ▶ **Adversarial training** in linear regression:

$$\max_{\|\Delta x_i\| \leq \delta} (y_i - \beta^T (x_i + \Delta x_i))^2$$

# Adversarially-trained linear regression

- ▶ **Linear regression:**

$$\min_{\beta} \sum_{i=1}^{\#train} (y_i - \beta^T x_i)^2$$

- ▶ **Adversarial training** in linear regression:

$$\min_{\beta} \sum_{i=1}^{\#train} \max_{\|\Delta x_i\| \leq \delta} (y_i - \beta^T (x_i + \Delta x_i))^2$$

## Adversarially-trained linear regression

$$\sum_{i=1}^{\#train} \max_{\|\Delta x_i\| \leq \delta} (y_i - (x_i + \Delta x_i)^T \beta)^2$$

## Adversarially-trained linear regression

$$\sum_{i=1}^{\#train} \max_{\|\Delta x_i\| \leq \delta} (y_i - (x_i + \Delta x_i)^T \beta)^2$$

It can be **rewritten** as:

$$\sum_{i=1}^{\#train} \left( |y_i - x_i^T \beta| + \delta \|\beta\|_* \right)^2$$

where  $\|\cdot\|_*$  is the **dual norm**.

## Adversarially-trained linear regression

$$\sum_{i=1}^{\#train} \max_{\|\Delta x_i\|_\infty \leq \delta} (y_i - (\mathbf{x}_i + \Delta \mathbf{x}_i)^\top \boldsymbol{\beta})^2$$

It can be **rewritten** as:

$$\sum_{i=1}^{\#train} \left( |y_i - \mathbf{x}_i^\top \boldsymbol{\beta}| + \delta \|\boldsymbol{\beta}\|_1 \right)^2$$

where  $\|\cdot\|_1$  is the **dual norm**.

## Similarities with Lasso

- ▶  $l_\infty$ -adversarial attacks:

$$\sum_{i=1}^{\#train} \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| + \delta \|\boldsymbol{\beta}\|_1 \right)^2$$

- ▶ Lasso:

$$\sum_{i=1}^{\#train} \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| \right)^2 + \lambda \|\boldsymbol{\beta}\|_1.$$

## Similarities with Lasso

- ▶  $l_\infty$ -adversarial attacks:

$$\sum_{i=1}^{\#train} \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| + \delta \|\boldsymbol{\beta}\|_1 \right)^2$$

- ▶ Lasso:

$$\sum_{i=1}^{\#train} \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| \right)^2 + \lambda \|\boldsymbol{\beta}\|_1.$$



## Similarities with Lasso

- ▶  $l_\infty$ -adversarial attacks:

$$\sum_{i=1}^{\#train} \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| + \delta \|\boldsymbol{\beta}\|_1 \right)^2$$

- ▶ Lasso:

$$\sum_{i=1}^{\#train} \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| \right)^2 + \lambda \|\boldsymbol{\beta}\|_1.$$

### Main results:

- #1. **Map**  $\lambda \leftrightarrow \delta$  for which they yield the **same result**.

## Similarities with Lasso

- ▶  $l_\infty$ -adversarial attacks:

$$\sum_{i=1}^{\#train} \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| + \delta \|\boldsymbol{\beta}\|_1 \right)^2$$

- ▶ Lasso:

$$\sum_{i=1}^{\#train} \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| \right)^2 + \lambda \|\boldsymbol{\beta}\|_1.$$

### Main results:

- #1. **Map**  $\lambda \leftrightarrow \delta$  for which they yield the **same result**.
- #2. **More parameters than data**: abrupt transition into interpolation.

## Similarities with Lasso

- ▶  $l_\infty$ -adversarial attacks:

$$\sum_{i=1}^{\#train} \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| + \delta \|\boldsymbol{\beta}\|_1 \right)^2$$

- ▶ Lasso:

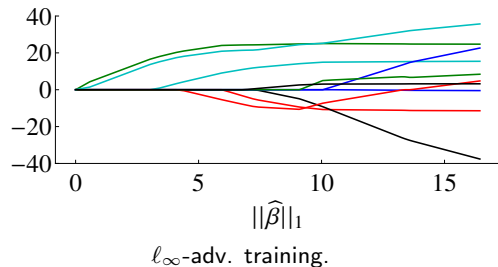
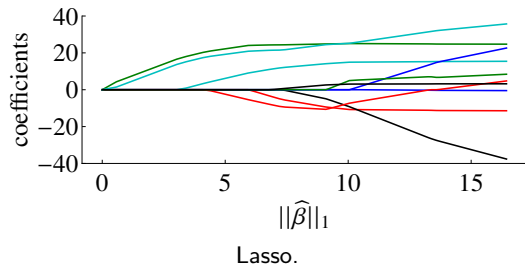
$$\sum_{i=1}^{\#train} \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| \right)^2 + \lambda \|\boldsymbol{\beta}\|_1.$$

### Main results:

- #1. **Map**  $\lambda \leftrightarrow \delta$  for which they yield the **same result**.
- #2. **More parameters than data**: abrupt transition into interpolation.
- #3. **Optimal choice** of  $\delta$  independent on noise level.

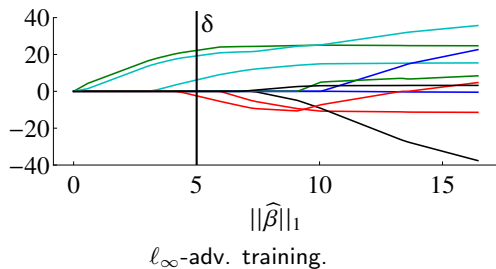
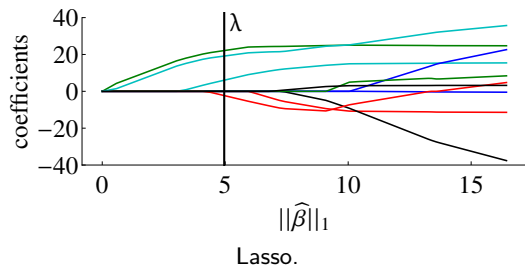
# # 1. Equivalence with Lasso

Map  $\lambda \leftrightarrow \delta$  for which they yield the **same result**.



## # 1. Equivalence with Lasso

Map  $\lambda \leftrightarrow \delta$  for which they yield the **same result**.



The that yield the **same result** are **not** necessarily the same, i.e.:  $\delta \neq \lambda$

## # 2. More parameters than data

**Lasso:** transition **only in the limit**

$$\lambda \rightarrow 0^+ \Rightarrow \text{Mean square error} \rightarrow 0$$

## # 2. More parameters than data

**Lasso:** transition **only in the limit**

$$\lambda \rightarrow 0^+ \Rightarrow \text{Mean square error} \rightarrow 0$$

**Adversarial training:**

$$\delta \in (0, \text{threshold}] \Rightarrow \text{Mean square error} = 0$$

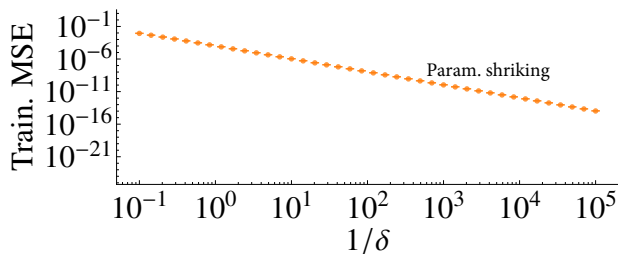
## # 2. More parameters than data

**Lasso:** transition **only in the limit**

$$\lambda \rightarrow 0^+ \Rightarrow \text{Mean square error} \rightarrow 0$$

**Adversarial training:**

$$\delta \in (0, \text{threshold}] \Rightarrow \text{Mean square error} = 0$$





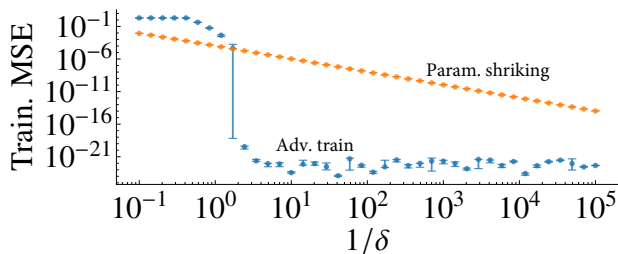
## # 2. More parameters than data

**Lasso:** transition **only in the limit**

$$\lambda \rightarrow 0^+ \Rightarrow \text{Mean square error} \rightarrow 0$$

**Adversarial training:**

$$\delta \in (0, \text{threshold}] \Rightarrow \text{Mean square error} = 0$$



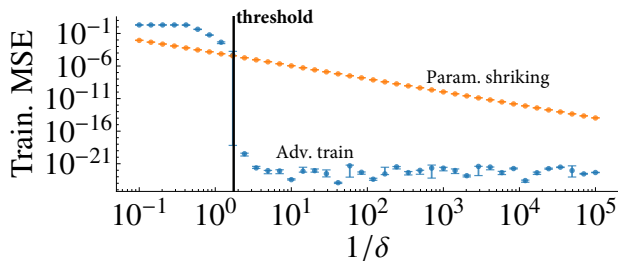
## # 2. More parameters than data

**Lasso:** transition **only in the limit**

$$\lambda \rightarrow 0^+ \Rightarrow \text{Mean square error} \rightarrow 0$$

**Adversarial training:**

$$\delta \in (0, \text{threshold}] \Rightarrow \text{Mean square error} = 0$$



## # 2. Equivalence with minimum norm interpolator

For  $\delta \in (0, \text{threshold}]$ , the minimum-norm interpolator is the solution to adversarial training.

## # 2. Equivalence with minimum norm interpolator

For  $\delta \in (0, \text{threshold}]$ , the minimum-norm interpolator is the solution to adversarial training.

### Relevance

Connect **adversarial training** with **double descent** and **benign overfitting**

## # 3. Invariance to noise levels

*To obtain near-oracle performance.*

▶ *Lasso:*

$$\lambda \propto \sigma \sqrt{\log(\#params)/\#train}$$

▶  *$\ell_\infty$ -adversarial attack:*

$$\delta \propto \sqrt{\log(\#params)/\#train}$$

### # 3. Invariance to noise levels

To obtain near-oracle performance.

► Lasso:

$$\lambda \propto \underbrace{\sigma}_{\text{unknown}} \sqrt{\log(\#params)/\#train}$$

►  $\ell_\infty$ -adversarial attack:

$$\delta \propto \sqrt{\log(\#params)/\#train}$$

#### Data model

$$y = \underbrace{x^T \beta^*}_{\text{signal}} + \underbrace{\sigma}_{\text{noise std.}} \varepsilon.$$

**arXiv:2310.10807**

▶  $\ell_2$ -adv. attacks and **ridge regression**.

**arXiv:2310.10807**

- ▶  $\ell_2$ -adv. attacks and **ridge regression**.
- ▶ Generalization to **other loss** functions



**arXiv:2310.10807**

- ▶  $\ell_2$ -adv. attacks and **ridge regression**.
- ▶ Generalization to **other loss** functions
- ▶ Connection to **robust regression** and  $\sqrt{\text{Lasso}}$ .