

# Adversarial training in linear regression

Antônio Horta Ribeiro  
Uppsala University

Seminar on Advances in Probabilistic Machine Learning  
Aalto University and ELLIS unit Helsinki  
Online, Nov 2022

# Outline

Motivation

Adversarial training

Robustness in high-dimensions

# Electrocardiogram exam

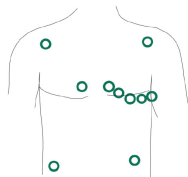
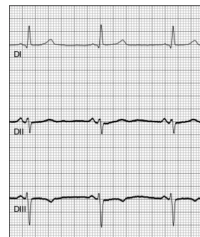
**Goal:** *Build data-driven ECG analysis tools.*

- ▶ The ECG is the major diagnostic tool.
- ▶ Cardiovascular diseases: 32% of all deaths (GBD 2019).
- ▶ Example. CODE dataset: annotated historical data  $n = 1.6\text{M}$  patients



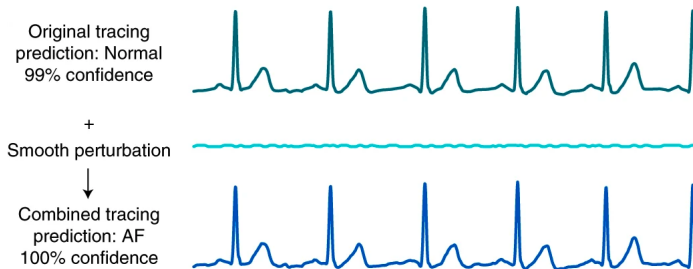
- ▶ Model for automatic diagnosis:

**A. H. Ribeiro**, M.H. Ribeiro, Paixão, G.M.M., et al. "Automatic diagnosis of the 12-lead ECG using a deep neural network," Nature Communications, 2020



**Left:** ECG signal **Right:** Electrode placement.

# Adversarial examples



**Figure:** Effect of adversarial examples on ECG Classification.

Source: Han, X., Hu, Y., Foschini, L. et al. Deep learning models for electrocardiograms are susceptible to adversarial attack. *Nature Medicine* 26, 360–363 (2020).

**Adversarial Training:** *Model is trained training on samples that have been modified by an adversary.*

## Adversarial training

*“Is it fundamentally different than other regularization methods?”*

Surprises in adversarially-trained linear regression (2022). **Antônio H. Ribeiro**, Dave Zachariah, Thomas B. Schön. arXiv:2205.12695.

## Adversarial robustness

*“what is the role of high-dimensionality in model robustness?”*

Overparameterized Linear Regression under Adversarial Attacks (2022). **Antônio H. Ribeiro**, Thomas B. Schön. arXiv:2204.06274

# Outline

Motivation

Adversarial training

Robustness in high-dimensions

# Framework: Linear regression

*Simplest case where adversarial vulnerability has been observed.*

I. J. Goodfellow, J. Shlens, C. Szegedy, "Explaining and Harnessing Adversarial Examples", ICLR 2015

D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Ma, "Robustness May Be At Odds with Accuracy," ICLR, p. 23, 2019.

- ▶ Training dataset:

$$(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n) \Rightarrow \hat{\beta}$$

- ▶ Model prediction

$$\hat{\mathbf{y}} = \hat{\beta}^T \mathbf{x}$$

- ▶ Error( $\hat{\beta}$ ) =  $|\mathbf{y} - \mathbf{x}^T \hat{\beta}|$

- ▶ Adv-error( $\hat{\beta}$ ) =  $\max_{\|\Delta \mathbf{x}\| \leq \delta} \left| \mathbf{y} - (\mathbf{x} + \Delta \mathbf{x})^T \hat{\beta} \right|$

# Adversarial training

Empirical risk minimization:

$$\min_{\beta} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2$$

Adversarial training:

$$\min_{\beta} \frac{1}{n} \sum_{i=1}^n \max_{\|\Delta \mathbf{x}_i\| \leq \delta} (y_i - (\mathbf{x}_i + \Delta \mathbf{x}_i)^T \beta)^2$$



# Adversarial error in linear regression

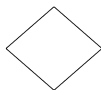
- ▶  $\text{Error}(\hat{\beta}) = |y - \mathbf{x}^T \hat{\beta}|$
- ▶  $\text{Adv-error}(\hat{\beta}) = \max_{\|\Delta \mathbf{x}\| \leq \delta} \left| y - (\mathbf{x} + \Delta \mathbf{x})^T \hat{\beta} \right|$
- ▶ *Dual formula for the adversarial error*

$$\left( \text{Adv-error}(\hat{\beta}) \right)^2 = \left( |\text{Error}(\hat{\beta})| + \delta \|\hat{\beta}\|_* \right)^2$$

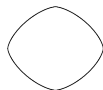
- ▶ where  $\|\cdot\|_*$  is the dual norm.

## $\ell_p$ -adversarial attacks

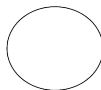
- ▶  $\ell_\infty$ -adversarial attack:  $\{\|\Delta x\|_\infty \leq \delta\} \Rightarrow$  dual norm:  $\|\Delta x\|_1$
- ▶  $\ell_2$ -adversarial attack:  $\{\|\Delta x\|_2 \leq \delta\} \Rightarrow$  dual norm:  $\|\Delta x\|_2$
- ▶  $\ell_p$ -adversarial attack:  $\{\|\Delta x\|_p \leq \delta\} \Rightarrow$  dual norm:  $\|\Delta x\|_q$   
for  $1/p + 1/q = 1$



$\ell_1$



$\ell_{1.5}$



$\ell_2$



$\ell_{20}$



$\ell_\infty$

# Consequences to adversarial training

- Adversarial training,

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

can be reformulated as

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

The above expression is **convex**

# Lasso and $\ell_\infty$ -adversarial training

- ▶  $\ell_\infty$ -adversarial training:

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

- ▶ Lasso:

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

# Ridge regression and $\ell_2$ -adversarial training

- ▶  $\ell_2$ -adversarial training:

$$\frac{1}{n} \sum_{i=1}^n \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| + \delta \|\boldsymbol{\beta}\|_2 \right)^2$$

- ▶ Ridge:

$$\frac{1}{n} \sum_{i=1}^n \left( |y_i - \mathbf{x}_i^T \boldsymbol{\beta}| \right)^2 + \delta \|\boldsymbol{\beta}\|_2^2$$

# Diabetes example

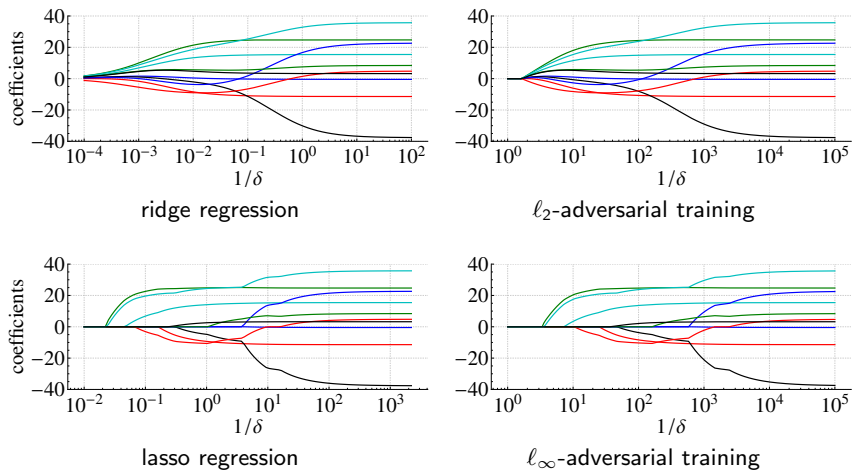


Figure: Regularization paths.

# Overparametrized models and interpolators

*Can a model perfectly fit the training data and still generalize well?*

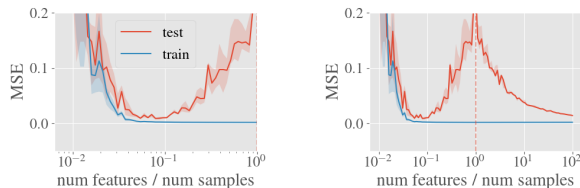
## ► Benign overfitting

P. L. Bartlett, P. M. Long, G. Lugosi, and A. Tsigler, "Benign overfitting in linear regression," *Proceedings of the National Academy of Sciences*, vol. 117, no. 48, pp. 30063–30070, Apr. 2020.

## ► Double descent

M. Belkin, D. Hsu, S. Ma, and S. Mandal, "Reconciling modern machine-learning practice and the classical bias–variance trade-off," *PNAS* (2019)

## ► Example:



**Figure:** nonlinear ARX mean squared error (MSE).

**A. H. Ribeiro**, J. N. Hendriks, A. G. Wills, T. B. Schön. "Beyond Occam's Razor in System Identification: Double-Descent when Modeling Dynamics". IFAC SYSID 2021. *Honorable mention: Young author award*

# Minimum-norm solution

## Minimum $\ell_2$ -norm solution

$$\min_{\beta} \|\beta\|_2 \quad \text{subject to} \quad X\beta = y$$

- ▶ Gradient descent in linear regression converges to  $\hat{\beta}^{\min-\ell_2}$ .
- ▶ Ridge  $\hat{\beta}^{\text{ridge}}(\delta) \rightarrow \hat{\beta}^{\min-\ell_2}$  as  $\delta \rightarrow 0^+$ .

## Minimum $\ell_1$ -norm solution

$$\min_{\beta} \|\beta\|_1 \quad \text{subject to} \quad X\beta = y$$

- ▶ Basis pursuit: i.e. allow you to recover sparse signals.
- ▶ Ridge  $\hat{\beta}^{\text{lasso}}(\delta) \rightarrow \hat{\beta}^{\min-\ell_1}$  as  $\delta \rightarrow 0^+$  (LARS algorithm)..



# Interpolation for finite $\delta$

## Theorem

For  $0 < \delta < \bar{\delta}$ , adversarial training is minimized at some  $\hat{\beta}$  that satisfies:

$$X\hat{\beta} = y$$

## Corollary

$\hat{\beta}^{\text{min-}\ell_2}$  is the solution to  $\ell_2$ -adversarial training for all  $0 < \delta < \bar{\delta}$ .

## Corollary

$\hat{\beta}^{\text{min-}\ell_1}$  is the solution to  $\ell_\infty$ -adversarial training for all  $0 < \delta < \bar{\delta}$ .

# Overparametrized model

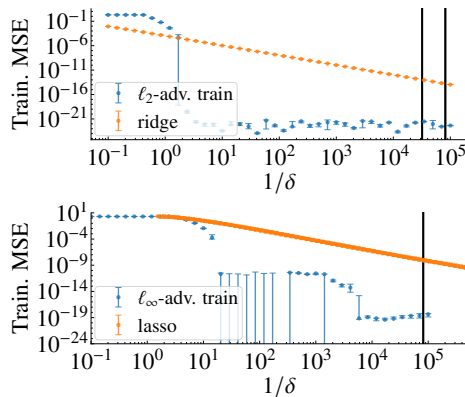


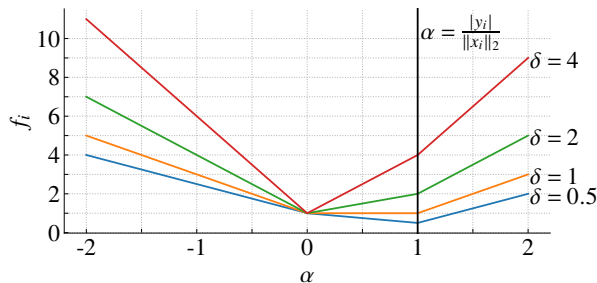
Figure: Training MSE vs regularization parameter.

# Discussion

- ▶ New interpretation for minimum-norm solution.
- ▶ Distinct behavior from other parameter shrinking methods (overparametrized).
- ▶ Explanation for abrupt transitions. Let:

$$f_i(\beta) = |y_i - \mathbf{x}_i^T \beta| + \delta \|\beta\|_2.$$

and assume  $|y_i| = \|\mathbf{x}_i\|_2 = 1$



# Outline

Motivation

Adversarial training

Robustness in high-dimensions

*Can a model perfectly fit the data and still be robust?*

# Analysing adversarial robustness

From:

$$\mathbb{E} \left[ \left( \text{Adv-error}(\hat{\beta}) \right)^2 \right] = \mathbb{E} \left[ \left( |\text{Error}(\hat{\beta})| + \delta \|\hat{\beta}\|_* \right)^2 \right]$$

It follows:

$$\mathbb{E}[\text{Error}(\hat{\beta})^2] + \delta^2 \|\hat{\beta}\|_*^2 \leq \mathbb{E}[(\text{Adv. error})^2] \leq \left( \sqrt{\mathbb{E}[\text{Error}(\hat{\beta})^2]} + \delta \|\hat{\beta}\|_* \right)^2.$$

# Application: plug-and-play from other analysis

Analysing minimum norm interpolation:

$$(\mathbf{x}_i, \epsilon_i) \sim P_{\mathbf{x}} \times P_{\epsilon}, \quad y_i = \mathbf{x}_i^{\top} \boldsymbol{\beta} + \epsilon_i,$$

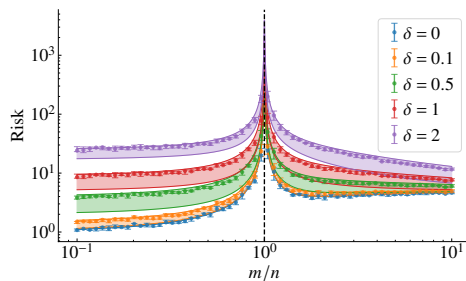


Figure: Adversarial risk vs number of features  $m$ .

T. Hastie, A. Montanari, S. Rosset, and R. J. Tibshirani, "Surprises in High-Dimensional Ridgeless Least Squares Interpolation," *Annals of Statistics*. 50(2): 949-986 (2022).

# Is robustness at odds with accuracy?

Can a good model to be arbitrarily vulnerable to adversarial attacks as you add more features?

## Proposition

If  $\mathbb{E}[\text{Error}(\hat{\beta})^2] < M$  as  $\# \text{features} \rightarrow \infty$ :

$$\mathbb{E}[(\text{Adv. error}(\hat{\beta}))^2] \rightarrow \infty$$

**if and only if**

$$\delta \|\hat{\beta}\|_* \rightarrow \infty.$$



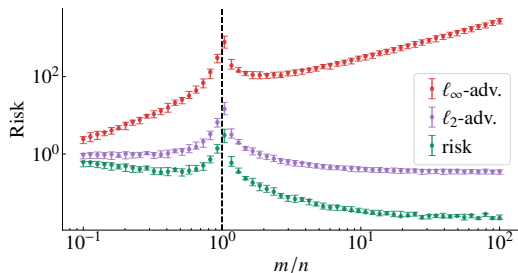
## Example

Minimum  $\ell_2$ -norm interpolator and Gaussian features:

$$\|\hat{\beta}\|_1 = \mathcal{O}(1) \quad \|\hat{\beta}\|_2 = \mathcal{O}(1/\sqrt{m})$$

Now, if we scale

$$\delta \propto \mathbb{E}\|\mathbf{x}\|_2 = \mathcal{O}(\sqrt{m}).$$



**Figure:** Adv. risk.

I. J. Goodfellow, J. Shlens, C. Szegedy, "Explaining and Harnessing Adversarial Examples", ICLR 2015

D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Ma, "Robustness May Be At Odds with Accuracy," ICLR, p. 23, 2019.

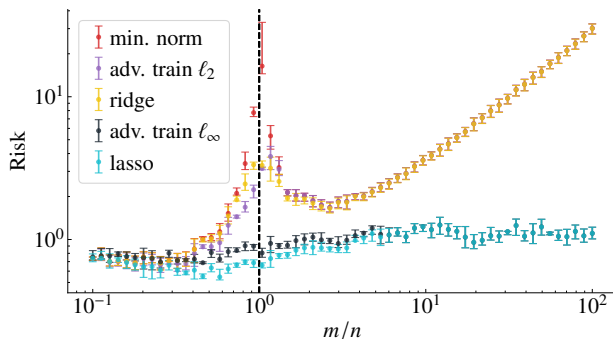
# The effect of regularization

- ▶ Ridge,  $\ell_2$ -adversarial training and min- $\ell_2$ -norm solution

$$\|\hat{\beta}\|_1 = \mathcal{O}(1)$$

- ▶ Lasso,  $\ell_\infty$ -adversarial training and min  $\ell_2$ -norm solution

$$\|\hat{\beta}\|_1 = \mathcal{O}(1/\sqrt{m})$$



**Figure:** Adversarial  $\ell_\infty$  risk and  $\delta \propto \mathbb{E}\|\mathbf{x}\|_2$ .

# Summary

- ▶ Dual formula for the adversarial error:

$$\left(\text{Adv-error}(\hat{\beta})\right)^2 = \left(|\text{Error}(\hat{\beta})| + \delta \|\hat{\beta}\|_*\right)^2$$

- ▶ Consequences to adversarial training:
  - ▶ Convex formula / Similarities with parameter shrink methods
  - ▶  $\ell_\infty$ -adversarial training  $\Rightarrow$  sparse solutions
  - ▶ Can interpolate for disturbance bounded by  $\delta > 0$ .
- ▶ Consequences to adversarial robustness
  - ▶ Simplify analysis of adversarial robustness.
  - ▶ Sufficient and necessary conditions for good models to be vulnerable to adversary .

**Thank you!**

✉ [antonio.horta.ribeiro@it.uu.se](mailto:antonio.horta.ribeiro@it.uu.se)

🌐 [antonior92.github.io](https://antonior92.github.io)