

# 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

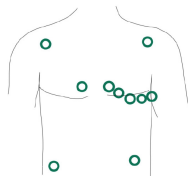
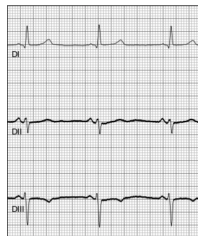
**Goal:** *Build data-driven ECG analysis tools.*

- | 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



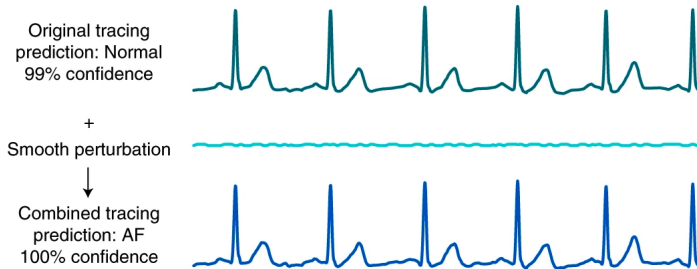
- | Model for automatic diagnosis:

**A. H. Ribeiro**, M.H. Ribeiro, Paixao, 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:

$$((x_1; y_1); (x_2; y_2); \dots; (x_n; y_n)) \quad b$$

| Model prediction

$$\hat{y} = b^T x$$

| Error( $b$ ) =  $\sum_j y_j - x^T b_j$

| Adv-error( $b$ ) =  $\max_k x_k \quad y_k - (x + \Delta x)^T b$

# Adversarial training

Empirical risk minimization:

$$\min \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T)^2$$

Adversarial training:

$$\min \frac{1}{n} \sum_{i=1}^n \max_k (y_i - (x_i + \Delta x_i)^T)^2$$



# Adversarial error in linear regression

- |  $\text{Error}(\mathbf{b}) = \sum_j y_j - \mathbf{x}_j^T \mathbf{b}$

- |  $\text{Adv-error}(\mathbf{b}) = \max_{\|\Delta \mathbf{x}\|_k} \sum_j y_j - (\mathbf{x}_j + \Delta \mathbf{x})^T \mathbf{b}$

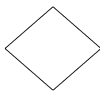
- | *Dual formula for the adversarial error*

$$\text{Adv-error}(\mathbf{b})^2 = \sum_j \text{Error}(\mathbf{b})_j^2 + \|\mathbf{b}\|_k^2$$

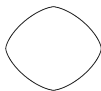
- | where  $\|\cdot\|_k$  is the dual norm.

# $\ell_p$ -adversarial attacks

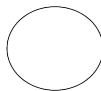
- |  $\ell_1$ -adversarial attack:  $f(x) \Delta x$   $k_1$        $g)$  dual norm:  $k \Delta x$   $k_1$
- |  $\ell_2$ -adversarial attack:  $f(x) \Delta x$   $k_2$        $g)$  dual norm:  $k \Delta x$   $k_2$
- |  $\ell_p$ -adversarial attack:  $f(x) \Delta x$   $k_p$        $g)$  dual norm:  $k \Delta x$   $k_q$   
for  $1/p + 1/q = 1$



$\ell_1$



$\ell_{1.5}$



$\ell_2$



$\ell_{20}$



$\ell_1$

# Consequences to adversarial training

| Adversarial training,

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

can be reformulated as

$$\frac{1}{n} \sum_{i=1}^n \max_{\| \Delta x \|_k \leq \epsilon} (y_i - x_i^T w - \Delta x^T w)^2$$

The above expression is **convex**

# Lasso and $\ell_1$ -adversarial training

|  $\ell_1$ -adversarial training:

$$\frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_1$$

| Lasso:

$$\frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_1$$

# Ridge regression and $\ell_2$ -adversarial training

|  $\ell_2$ -adversarial training:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{j})^2 + k \|\mathbf{j}\|_2^2$$

| Ridge:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{j})^2 + k \|\mathbf{j}\|_2^2$$

# Diabetes example

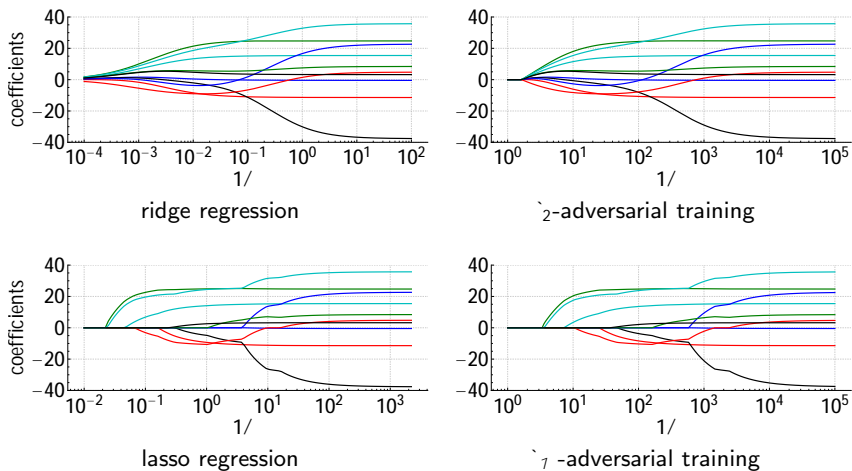


Figure: Regularization paths.

# Overparametrized models and interpolators

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

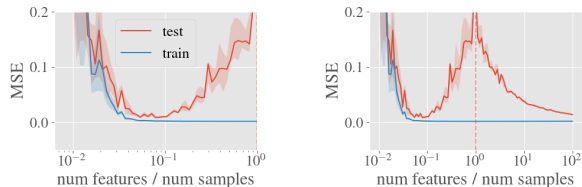
## Benign overfitting

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

$$\min_k \|k\|_2 \quad \text{subject to} \quad Xk = y$$

- | Gradient descent in linear regression converges to  $b^{\min \ell_2}$ .
- | Ridge  $b^{\text{ridge}}(\lambda)$   $\rightarrow$   $b^{\min \ell_2}$  as  $\lambda \rightarrow 0^+$ .

## Minimum $\ell_1$ -norm solution

$$\min_k \|k\|_1 \quad \text{subject to} \quad Xk = y$$

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



# Interpolation for finite

## Theorem

For  $0 < \epsilon < \bar{\epsilon}$ , adversarial training is minimized at some  $b$  that satisfies:

$$x^b = y$$

## Corollary

$b_{\min-\ell_2}$  is the solution to  $\ell_2$ -adversarial training for all  $0 < \epsilon < \bar{\epsilon}$ .

## Corollary

$b_{\min-\ell_1}$  is the solution to  $\ell_1$ -adversarial training for all  $0 < \epsilon < \bar{\epsilon}$ .

# Overparametrized model

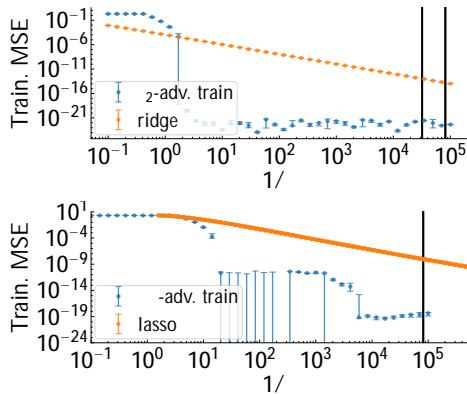


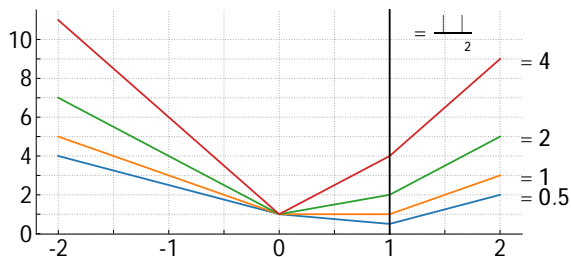
Figure: Training MSE vs regularization parameter.

# Discussion

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

$$f_i(\beta) = |y_i - x_i^T \beta| + k_1 \|\beta\| + k_2 \|\beta\|_2$$

and assume  $|y_i| = k_1 k_2 = 1$



# Outline

Motivation

Adversarial training

Robustness in high-dimensions

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

# Analysing adversarial robustness

From:

$$\mathbb{E} \|\text{Adv-error}(\mathbf{b})\|^2 = \mathbb{E} \|\text{Error}(\mathbf{b})\|^2 + \|\mathbf{b}\|^2$$

It follows:

$$\mathbb{E}[\|\text{Error}(\mathbf{b})\|^2] + \|\mathbf{b}\|^2 = \mathbb{E}[\|\text{Adv. error}\|^2] \iff \frac{\mathbb{E}[\|\text{Adv. error}\|^2]}{\mathbb{E}[\|\text{Error}(\mathbf{b})\|^2] + \|\mathbf{b}\|^2} :$$

# Application: plug-and-play from other analysis

Analysing minimum norm interpolation:

$$(x_i; i) \quad P_X \quad P; \quad y_i = x_i^T + i;$$

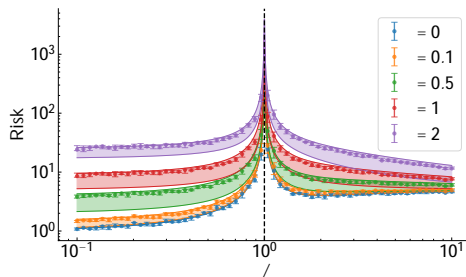


Figure: Adversarial risk vs number of features  $m$ .

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 be arbitrarily vulnerable to adversarial attacks as you add more features?

## Proposition

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

$$E[(\text{Adv. error}(\mathbf{b}))^2] \rightarrow \infty$$

**if and only if**

$$\|\mathbf{b}\|_k \rightarrow \infty$$



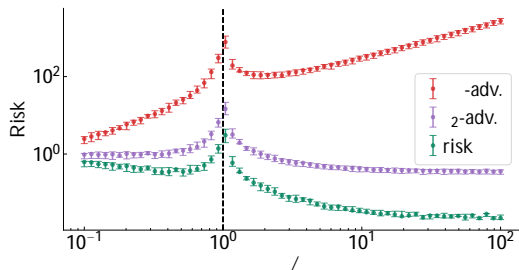
## Example

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

$$\|b\|_{k_1} = O(1) \quad \|b\|_{k_2} = O(1/\sqrt{m})$$

Now, if we scale

$$\|E\|_{k_2} = O(1/\sqrt{m}):$$



**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.

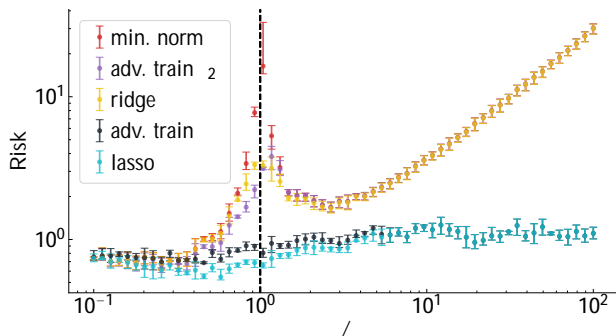
## The effect of regularization

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

$$\|k\|_{k_1} = O(1)$$

- Lasso,  $\ell_1$ -adversarial training and min  $\ell_2$ -norm solution

$$\|k\|_{k_1} = O(1 = \frac{\rho}{m})$$



**Figure:** Adversarial  $\ell_1$  risk and  $\|k\|_{k_2}$ .

# Summary

- | Dual formula for the adversarial error:

$$\text{Adv-error}(\mathbf{b})^2 = \|\text{Error}(\mathbf{b})\|^2 + \|\mathbf{b}\|^2$$

- | Consequences to adversarial training:

- | Convex formula / Similarities with parameter shrink methods
- |  $\ell_\infty$ -adversarial training ) sparse solutions
- | Can interpolate for disturbance bounded by  $\epsilon > 0$ .

- | 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

🌐 antonior92.github.io