

UPPSALA UNIVERSITET



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

Adversarial training in linear models

Adversarially-trained linear regression minimizes:

$$\sum_{i=1}^{n} \max_{\|\Delta x_i\| \le \delta} (y_i - (\mathbf{x}_i + \Delta x_i)^{\mathsf{T}} \boldsymbol{\beta})^2$$

Cost function is **convex**. It can be **rewritten** as: [2]

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

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

[2] A H Ribeiro and T B Schön, "Overparameterized Linear Regression under Adversarial Attacks", IEEE Transactions on Signal Processing, 2023

atta alra
attaalra
attacks
Cost function
$\sum_{i=1}^{n} \left(y_i - \boldsymbol{x}_i^{T}\boldsymbol{\beta} + \delta \boldsymbol{\beta} _2 \right)^2$
l attacks
Cost function
$\sum_{i=1}^{n} \left(y_i - \boldsymbol{x}_i^{T}\boldsymbol{\beta} + \delta \ \boldsymbol{\beta}\ _1 \right)^2$
-

Regularization properties of adversarially-trained linear regression

Antônio H. Ribeiro^{*,1}, Dave Zachariah¹, Francis Bach², Thomas B. Schön¹ ¹Uppsala University (Sweden), ²ÍNRIA/PSL Research University (France)

* contact: antonio.horta.ribeiro@it.uu.se

Motivation

How does adversarial training compare with other regularization methods?

Regularization methods:

- . Parameter shrinking methods.
- 2. Minimum-norm interpolators.
- $\frac{1}{2}$ 3. Robust regression and $\sqrt{\text{Lasso}}$.

#1: Equivalence with parameter shrinking

Backgroud:

- $\min_{\boldsymbol{\beta}} \sum_{i=1}^{n} \left(|\boldsymbol{y}_i \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{\beta}| \right)^2 + \lambda \|\boldsymbol{\beta}\|_2^2.$ • Ridge:
- $\min_{\boldsymbol{\beta}} \sum_{i=1}^{n} \left(|y_i \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{\beta}| \right)^2 + \lambda \|\boldsymbol{\beta}\|_1.$ • Lasso:

Our result: map $\lambda \leftrightarrow \delta$ for which the results are asymptotically equiva**lent** if $\mathbb{E}[x] = 0$ and $x \sim -x$.

Numerical experiments:









• Used in the study of **double-descent** [3] and **benign-overfitting**.





Invariance to noise magnitude

Background: If data is generated as

$$y = \underbrace{x^{\top} \beta^{*}}_{\text{signal}} + \sigma \underbrace{\varepsilon}_{\mathcal{N}(0,1)}$$

For Lasso, **near-oracle** performance is attained with:

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

Our result: For ℓ_{∞} -adv. attack, **near-oracle** performance is attained with

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

That is, the adv. radius δ can be set without knowledge of the noise magnitude.

#3: Relation to robust regression and \sqrt{Lasso}

 $\sqrt{\text{Lasso:}}$ Has the same property. It is the estimator: [4]

$$\min_{\boldsymbol{\beta}} \sqrt{\sum_{i=1}^{n} |y_i - \boldsymbol{x}_i^{\top} \boldsymbol{\beta}|^2} + \lambda \|\boldsymbol{\beta}\|_1.$$

and it attains **near-oracle** performance if:

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

Robust regression: [5]

$$\min_{\beta} \max_{\Delta \in \mathcal{S}} \|y - (X + \Delta)\beta\|_2.$$

It is equivalent to $\sqrt{\text{Lasso}}$ if columns are constrained: $\|\Delta_{i,:}\|_2 \leq \delta, \forall i$.

Our result: Robust regression is **equivalent** to adversarial training if **rows** are constrained: $\|\Delta_{:,j}\|_2 \leq \delta, \forall j$.

[4] A. Belloni, V. Chernozhukov, and L. Wang, "Square-Root Lasso: Pivotal Recovery of Sparse Signals via Conic Programming," Biometrika 2011. [5] H. Xu, C. Caramanis, and S. Mannor, "Robust regression and lasso," NIPS 2008

Future work

- Error-in-variables: adv. training making it more robust to input disturbance Δx that are stochastic.
- **Tailored solver:** enableuse in high-dimensional applications (genetics).
- **Nonlinear models:** most results still hold for inputs in infinite spaces.