

Overparametrized Linear Regression under Adversarial Attacks

Antônio Horta Ribeiro

Uppsala University (Sweden)

SIERRA group meeting

INRIA Paris, March 2023

Outline

Motivation

Robustness in high-dimensions

Adversarial training

The state of Minas Gerais and telehealth

| Minas Gerais



The state of Minas Gerais and telehealth

- | Minas Gerais
 - | approximately the same area as France.



The state of Minas Gerais and telehealth

- | Minas Gerais
 - | approximately the same area as France.
 - | 853 municipalities



The state of Minas Gerais and telehealth

- | Minas Gerais
 - | approximately the same area as France.
 - | 853 municipalities
- | Telehealth center of Minas Gerais



The state of Minas Gerais and telehealth

- | Minas Gerais
 - | approximately the same area as France.
 - | 853 municipalities
- | Telehealth center of Minas Gerais
- | More than 4000 ECGs per day



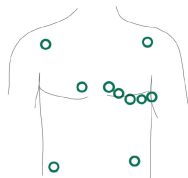
Electrocardiogram exam and machine learning

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

Electrocardiogram exam and machine learning

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

- | Cardiovascular diseases: 32% of all deaths (GBD 2019).

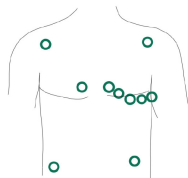


Left: ECG signal **Right:** Electrode placement.

Electrocardiogram exam and machine learning

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

- | Cardiovascular diseases: 32% of all deaths (GBD 2019).
- | The ECG is the major diagnostic tool.

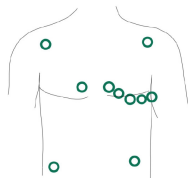


Left: ECG signal **Right:** Electrode placement.

Electrocardiogram exam and machine learning

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

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

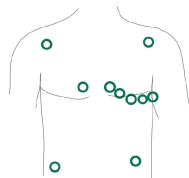
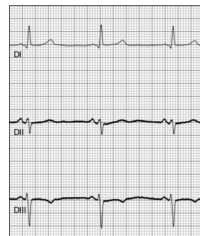
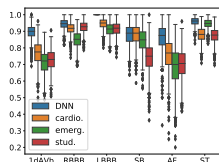


Left: ECG signal **Right:** Electrode placement.

Electrocardiogram exam and machine learning

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

- | Cardiovascular diseases: 32% of all deaths (GBD 2019).
- | The ECG is the major diagnostic tool.
- | CODE dataset: annotated historical data $n = 1.6M$ patients.
- | Model for automatic diagnosis:

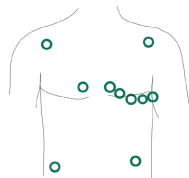
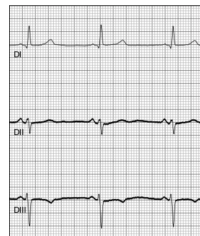
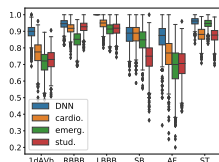


Left: ECG signal **Right:** Electrode placement.

Electrocardiogram exam and machine learning

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

- | Cardiovascular diseases: 32% of all deaths (GBD 2019).
- | The ECG is the major diagnostic tool.
- | CODE dataset: annotated historical data $n = 1.6M$ patients.
- | Model for automatic diagnosis:

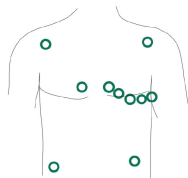
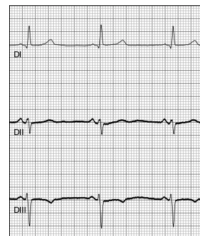
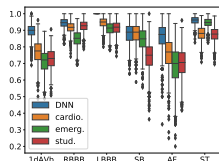


Left: ECG signal **Right:** Electrode placement.

Electrocardiogram exam and machine learning

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

- | Cardiovascular diseases: 32% of all deaths (GBD 2019).
- | The ECG is the major diagnostic tool.
- | CODE dataset: annotated historical data $n = 1.6M$ patients.
- | Model for automatic diagnosis:



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

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

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 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 robustness

"What is the role of high-dimensionality in model robustness?"

Overparameterized Linear Regression under Adversarial Attacks (2023). **Antônio H. Ribeiro**, Thomas B. Schön. *IEEE Transactions on Signal Processing* (preprint: arxiv.org/abs/2204.06274).

Adversarial robustness

"What is the role of high-dimensionality in model robustness?"

Overparameterized Linear Regression under Adversarial Attacks (2023). **Antônio H. Ribeiro**, Thomas B. Schön. *IEEE Transactions on Signal Processing* (preprint: arxiv.org/abs/2204.06274).

Adversarial training

"How is it connected to other regularization methods?"

Surprises in adversarially-trained linear regression (2022). **Antônio H. Ribeiro**, Dave Zachariah, Francis Bach, Thomas B. Schön. *Work in progress*.

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:

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

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

$$y = b^T x$$

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:

$$(\mathbf{x}_1; \mathbf{y}_1); (\mathbf{x}_2; \mathbf{y}_2); \dots; (\mathbf{x}_n; \mathbf{y}_n) \quad \mathbf{b}$$

| Model prediction

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

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

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

$$b = b^T x$$

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

| Adv-error(b) = $\max_{x_k} \sum_k y_k - (x + \delta x)^T b$

Adversarial error in linear regression

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

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

- | Dual formula for the adversarial error

$$\text{Adv-error}(\mathbf{b})^2 = |y - \mathbf{x}^T \mathbf{b}|^2 + \|\mathbf{b}\|_k^2$$

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

ℓ_p -adversarial attacks

f) ℓ_1 -adversarial attack: $\|x\|_1$ g) dual norm: $\|x\|_1$

ℓ_p -adversarial attacks

- | ℓ_1 -adversarial attack: $f(x) - f(x_0)$ g) dual norm: $\|x - x_0\|_1$
- | ℓ_2 -adversarial attack: $f(x) - f(x_0)$ g) dual norm: $\|x - x_0\|_2$

ℓ_p -adversarial attacks

- | ℓ_1 -adversarial attack: $f(x) = \|x\|_1$ g) dual norm: $\|x\|_\infty$
- | ℓ_2 -adversarial attack: $f(x) = \|x\|_2$ g) dual norm: $\|x\|_2$
- | ℓ_p -adversarial attack: $f(x) = \|x\|_p$ g) dual norm: $\|x\|_q$
for $\frac{1}{p} + \frac{1}{q} = 1$

ℓ_1

$\ell_{1.5}$

ℓ_2

ℓ_{20}

ℓ_∞

Outline

Motivation

Robustness in high-dimensions

Adversarial training

What is the role of high-dimensionality in model robustness?

I High-dimensionality as a source of vulnerability:

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, 2019.

J. Gilmer et al., \Adversarial Spheres," arXiv:1801.02774, Sep. 2018.

What is the role of high-dimensionality in model robustness?

I High-dimensionality as a source of vulnerability:

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, 2019.

J. Gilmer et al., \Adversarial Spheres," arXiv:1801.02774, Sep. 2018.

I High-dimensionality as a source of robustness:

S. Bubeck and M. Sellke, \A Universal Law of Robustness via Isoperimetry," Advances in Neural Information Processing Systems, 2021

How neural networks can perform well?

Language models

Image classification

Figure: Models number of parameters

Sources: J. Simon (2021) "Large Language Models: A New Moore's Law?". Online (accessed: 2021-11-09). URL: huggingface.co/blog/large-language-models .

M. Tan and Q. V. Le (2019) "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," ICML

Overparametrized models

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

I 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.

Overparametrized models

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

I 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.

I 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)

Overparametrized models

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

I 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.

I 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)

I 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

Setup: minimum-norm interpolator

$$\min \|x\|_2 \quad \text{subject to} \quad X = y$$

- | Gradient descent in linear regression converges to b_{\min}^2 .

Setup: minimum-norm interpolator

$$\min \|k\|_{k_2} \quad \text{subject to} \quad X = y$$

- | Gradient descent in linear regression converges to b^{\min} .
- | Used to study benign overfitting in:

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.

Setup: minimum-norm interpolator

$$\min \|k\|_{k_2} \quad \text{subject to} \quad X = y$$

| Gradient descent in linear regression converges to b_{\min}^2 .

| Used to study benign overfitting in:

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.

| Use to study double descent in:

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

Analysing adversarial robustness

From:

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

Analysing adversarial robustness

From:

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

It follows that:

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

Double-descent in the adversarial loss

$$E[\ell_2\text{-adv. error}(\mathbf{b})^2] / E[\text{Error}(\mathbf{b})^2] + \|\mathbf{b}\|_2^2:$$

Double-descent in the adversarial loss

$$E[(\ell_2\text{-adv. error}(\mathbf{b}))^2] / E[\text{Error}(\mathbf{b})^2] + \|\mathbf{b}\|_2^2:$$

$\|\mathbf{b}\|_2$ also present a double descent behavior. As we increase the problem dimension, it becomes possible to find solutions with smaller norm.

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

Asymptotic results

Analysing minimum-norm interpolation:

$$(x_i; y_i) \sim P_x \quad P; \quad y_i = x_i^T \beta + \epsilon_i;$$

Figure: Adversarial risks 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).

High-dimensionality as a source of brittleness

| Model:

$$y \sim N(0; 1)$$

$$x_i \sim N(y; 1)$$

High-dimensionality as a source of brittleness

| Model:

$$y \sim N(0; 1)$$

$$x_i \sim N(y; 1)$$

| Optimal predictor: $b = \frac{h}{\# \text{ features}}; \dots; \frac{i}{\# \text{ features}}$

High-dimensionality as a source of brittleness

| Model:

$$y \sim N(0; 1)$$

$$x_i \sim N(y; 1)$$

| Optimal predictor: $b = \frac{1}{\# \text{ features}}; \quad ; \frac{1}{\# \text{ features}}$

| $= E \|x\|_2 / \sqrt{\# \text{ features}}$

High-dimensionality as a source of brittleness

| Model:

$$y \sim N(0; 1)$$

$$x_i \sim N(y; 1)$$

| Optimal predictor: $b = \frac{1}{\# \text{ features}}; \quad ; \frac{1}{\# \text{ features}}$

| $= E \|x\|_2 / \sqrt{\# \text{ features}}$

High-dimensionality as a source of brittleness

| Model:

$$y \sim \mathcal{N}(0; 1)$$

$$x_i \sim \mathcal{N}(y; 1)$$

| Optimal predictor: $b = \frac{h}{\# \text{ features}}; \quad ; \frac{1}{\# \text{ features}}$

| $= \mathbb{E} \|x\|_2 / \sqrt{\# \text{ features}}$

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.

Is robustness at odds with accuracy?

When good model become more vulnerable to adversarial attacks as we add features?

Proposition

If $E[\text{Error}(\mathbf{b})^2] < \epsilon$:

$$E[(\text{Adv. error}(\mathbf{b}))^2] \leq 1 \quad \text{as } \#\text{features} \rightarrow 1$$

if and only if

$$\|\mathbf{b}\|_2 \leq 1 :$$

Example: optimal predictor vulnerable to adversarial attacks

| Optimal predictor: $\mathbf{b} = \frac{1}{\sqrt{\text{\# features}}}$; \mathbf{i}

| $\mathbf{x} = \frac{\mathbf{e}_k}{\sqrt{\text{\# features}}}$

| For our example,

$$\|\mathbf{b}\|_1 = \sqrt{\text{\# features}}$$

hence

$$E[\text{\# adv. error}(\hat{\mathbf{b}})] = O(\sqrt{\text{\# features}})$$

Example: optimal predictor vulnerable to adversarial attacks

| Optimal predictor: $b = \frac{1}{\sqrt{\# \text{ features}}}$; i

| $x = \mathbb{E} \|x\|_2 / \sqrt{\# \text{ features}}$

| For our example,

$$\|b\|_1 = \sqrt{\# \text{ features}}$$

hence

$$\mathbb{E}[(\|x\|_1 - \text{adv. error}(b))^2] = O(\# \text{ features})$$

Mismatched example

| $\|x\|_1$ -adv. attack $\|x\|_1$.

| $x / \mathbb{E} \|x\|_2$.

:

Example: minimum ℓ_2 -norm interpolator

Minimum ℓ_2 -norm interpolator and Gaussian features:

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

Now, if we scale

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

Figure: Adv. risk.

The effect of adversarial training and regularization

Empirical risk minimization:

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

Adversarial training:

$$\min \frac{1}{n} \sum_{i=1}^n \max_{\|\mathbf{x}_i\|_k} (y_i - (\mathbf{x}_i + \mathbf{x}_i)^\top \mathbf{w})^2$$

The effect of regularization

- | Ridge and ℓ_2 -adversarial training

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

- | Lasso, ℓ_1 -adversarial training

$$\|b\|_{k_1} = O(\sqrt{\frac{p}{m}})$$

The effect of regularization

- | Ridge and ℓ_2 -adversarial training

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

- | Lasso, ℓ_1 -adversarial training

$$\|b\|_{k_1} = O(\sqrt{\frac{p}{m}})$$

Figure: Adversarial ℓ_1 risk and $\|E\|_{k_2}$.

Outline

Motivation

Robustness in high-dimensions

Adversarial training

Adversarial training in linear models

I Adversarial training,

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

Adversarial training in linear models

Adversarial training,

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

can be reformulated as

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

Lasso and ℓ_1 -adversarial training

| ℓ_1 -adversarial training:

$$\frac{1}{n} \sum_{i=1}^n |y_i - \mathbf{x}_i^T \mathbf{j}| + k \|\mathbf{k}_1\|^2$$

| Lasso:

$$\frac{1}{n} \sum_{i=1}^n |y_i - \mathbf{x}_i^T \mathbf{b}_j|^2 + k \|\mathbf{k}_1\|$$

Ridge regression and ℓ_2 -adversarial training

| ℓ_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

ridge regression

ℓ_2 -adversarial training

lasso regression

ℓ_1 -adversarial training

Figure: Regularization paths.

Minimum-norm solution

Minimum ℓ_2 -norm solution

$$\min \|x\|_2 \quad \text{subject to} \quad Xx = y$$

| Gradient descent in linear regression converges to $b_{\min}^{\ell_2}$.

| Used to study benign overfitting in:

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.

| Use to study double descent in:

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

Minimum-norm solution

Minimum ℓ_2 -norm solution

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

| Gradient descent in linear regression converges to $b_{\min}^{\ell_2}$.

| Used to study benign overfitting in:

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.

| Use to study double descent in:

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

Minimum ℓ_1 -norm solution

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

| Basis pursuit: i.e. allow you to recover sparse signals.

| Can also be used to study benign overfitting and double descent, i.e.,

F. Koehler, L. Zhou, D. J. Sutherland, and N. Srebro, "Uniform Convergence of Interpolators: Gaussian Width, Norm Bounds and Benign Overfitting," presented at the Advances in Neural Information Processing Systems, 2021

Minimum-norm interpolator and adversarial training

Theorem

Adversarial training is minimized at the minimum norm interpolator

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

$$i \quad 0 < \epsilon < \dots$$

Surprises in adversarially-trained linear regression (2022) António H. Ribeiro, Dave Zachariah, Francis Bach, Thomas B. Schön. Work in progress.

ℓ_2 -adversarial training vs ridge regression

Relation to min-norm solution

- | (corollary) $\mathbf{b}^{\min-\ell_2}$ is the solution to ℓ_2 -adversarial training if $0 < \lambda < \infty$.
- | Ridge $\mathbf{b}^{\text{ridge}}(\lambda) \rightarrow \mathbf{b}^{\min-\ell_2}$ as $\lambda \rightarrow 0^+$.

Figure: Training MSE vs regularization parameter.

ℓ_1 -adversarial trainings Lasso

Relation to min-norm solution

- | (corollary) $\mathbf{b}^{\min-\ell_1}$ is the solution to ℓ_1 -adversarial training if $0 < \lambda < \lambda_{\text{max}}$.
- | Lasso $\mathbf{b}^{\text{lasso}}(\lambda)$ \rightarrow $\mathbf{b}^{\min-\ell_1}$ as $\lambda \rightarrow 0^+$ (LARS algorithm)..

Figure: Training MSE vs regularization parameter.

Discussion

- | Distinct behavior from other parameter shrinking methods (overparametrized).

Discussion

- | Distinct behavior from other parameter shrinking methods (overparametrized).
- | Explanation for abrupt transitions. Let:

$$f_i(\theta) = |y_i - x_i^T \theta| + \lambda \|\theta\|_2$$

and assume $\|y_i\| = \|x_i\|_2 = 1$

New interpretation for minimum-norm interpolator

Figure: Threshold vs number of features m .

New interpretation for minimum-norm interpolator

Figure: Threshold vs number of features m .

- | Setup: $y_i = x_i^T \beta + \epsilon_i$, $x_i \sim N(0; r^2 I_m)$ and $\epsilon_i \sim N(0; \sigma^2)$
- | Increasing the number of features for minimum-norm interpolators increases the maximum disturbance in the corresponding adversarial training problems.

Summary

- | Dual formula for the adversarial error:

$$\left(\text{Adv-error}(\hat{\gamma})\right)^2 = \left(|\text{Error}(\hat{\gamma})| + \|\hat{\gamma}\|\right)^2$$

- | Consequences to adversarial robustness

- | Simplify analysis of adversarial robustness:

$$\mathbb{E} \left[\left(\text{Adv-error}(\hat{\gamma})\right)^2 \right] \propto \mathbb{E} \left[\text{Error}(\hat{\gamma})^2 \right] + \|\hat{\gamma}\|^2$$

- | Double descent can be observed in adversarial scenarios.

- | Sufficient and necessary conditions for good models to be vulnerable to adversary.

- | Consequences to adversarial training:

- | Convex formula / Similarities with parameter shrink methods

- | ℓ_1 -adversarial training \Rightarrow sparse solutions

- | Can interpolate for disturbance bounded by $\epsilon > 0$.

Thank you!

📍 INRIA de Paris - Room C407 (from now to mid-June)

✉ antonio.horta.ribeiro@it.uu.se

🌐 antonior92.github.io