

Overparametrized regression under l_2 adversarial attacks

Antônio H. Ribeiro, Thomas B. Schön

Uppsala University, Sweden



UPPSALA
UNIVERSITET

Workshop on the Theory of Overparameterized Machine Learning
TOPML, 2021

Overparametrized models can generalize effectively when train and test come from the same distribution...

Overparametrized models can generalize effectively when train and test come from the same distribution...

Can it also generalize effectively when there is a distribution shift?

Adversarial attacks

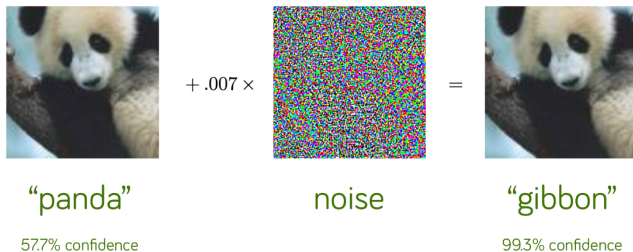


Figure: Illustration of adversarial attack. From: I.J. Goodfellow, J. Shlens, C.Szegedy , *“Explaining and Harnessing Adversarial Examples”*, ICLR 2015.

Linear regression under adversarial attacks

Model: Linear model

$$y = \mathbf{x}^T \boldsymbol{\beta} + \epsilon$$

Linear regression under adversarial attacks

Model: Linear model

$$y = \mathbf{x}^T \boldsymbol{\beta} + \epsilon$$

Estimated parameter: using train dataset (\mathbf{x}_i, y_i) , $i = 1, \dots, n$:

- ▶ Underparametrized:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \sum_i (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2$$

Linear regression under adversarial attacks

Model: Linear model

$$y = \mathbf{x}^T \boldsymbol{\beta} + \epsilon$$

Estimated parameter: using train dataset (\mathbf{x}_i, y_i) , $i = 1, \dots, n$:

- ▶ Underparametrized:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \sum_i (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2$$

- ▶ Overparametrized:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\boldsymbol{\beta}\|_2^2$$

subject to $y_i = \mathbf{x}_i^T \boldsymbol{\beta}$

for every i

Linear regression under adversarial attacks

Given a data point not seen during training (\mathbf{x}, y) .

Standard risk:

$$(y - \mathbf{x}^T \hat{\beta})^2$$

Adversarial risk:

$$(y - (\mathbf{x} + \Delta \mathbf{x})^T \hat{\beta})^2$$

Linear regression under adversarial attacks

Given a data point not seen during training (\mathbf{x}, y) .

Standard risk:

$$(y - \mathbf{x}^T \hat{\beta})^2$$

Adversarial risk:

$$\max_{\|\Delta \mathbf{x}\|_p \leq \delta} (y - (\mathbf{x} + \Delta \mathbf{x})^T \hat{\beta})^2$$

$\Delta \mathbf{x} \rightsquigarrow$ Adversarially generated disturbance

Linear regression under adversarial attacks

Given a data point not seen during training (\mathbf{x}, y) .

Standard risk:

$$R = E \left\{ (y - \mathbf{x}^T \hat{\beta})^2 \right\}$$

Adversarial risk:

$$R^{\text{adv}} = E \left\{ \max_{\|\Delta \mathbf{x}\|_p \leq \delta} (y - (\mathbf{x} + \Delta \mathbf{x})^T \hat{\beta})^2 \right\}$$

$\Delta \mathbf{x} \rightsquigarrow$ Adversarially generated disturbance

Linear regression under adversarial attacks

Given a data point not seen during training (\mathbf{x}, y) .

Standard risk:

$$R = E \left\{ (y - \mathbf{x}^T \hat{\beta})^2 \mid \underbrace{x_i, i = 1, \dots, n}_{\text{training inputs}} \right\}$$

Adversarial risk:

$$R^{\text{adv}} = E \left\{ \max_{\|\Delta \mathbf{x}\|_p \leq \delta} (y - (\mathbf{x} + \Delta \mathbf{x})^T \hat{\beta})^2 \mid \underbrace{x_i, i = 1, \dots, n}_{\text{training inputs}} \right\}$$

$\Delta \mathbf{x} \rightsquigarrow$ Adversarially generated disturbance

Bounds on the adversarial risk:

$$R + \delta^2 N_q + \sigma^2 \leq R^{\text{adv}} \leq \left(\sqrt{R} + \delta \sqrt{N_q} \right)^2 + \sigma^2.$$

Bounds on the adversarial risk:

$$R + \delta^2 N_q + \sigma^2 \leq R^{\text{adv}} \leq \left(\sqrt{R} + \delta \sqrt{N_q} \right)^2 + \sigma^2.$$

R^{adv} \rightsquigarrow Adversarial risk

Bounds on the adversarial risk:

$$R + \delta^2 N_q + \sigma^2 \leq R^{\text{adv}} \leq \left(\sqrt{R} + \delta \sqrt{N_q} \right)^2 + \sigma^2.$$

R^{adv} \rightsquigarrow Adversarial risk

R \rightsquigarrow Risk

Bounds on the adversarial risk:

$$R + \delta^2 N_q + \sigma^2 \leq R^{\text{adv}} \leq \left(\sqrt{R} + \delta \sqrt{N_q} \right)^2 + \sigma^2.$$

R^{adv} \rightsquigarrow Adversarial risk

R \rightsquigarrow Risk

$N_q = E\{\|\hat{\beta}\|_q^2\}$

Bounds on the adversarial risk:

$$R + \delta^2 N_q + \sigma^2 \leq R^{\text{adv}} \leq \left(\sqrt{R} + \delta \sqrt{N_q} \right)^2 + \sigma^2.$$

R^{adv} \rightsquigarrow Adversarial risk

R \rightsquigarrow Risk

$$N_q = E\{\|\hat{\beta}\|_q^2\}$$

\rightarrow For an ℓ_p attack $\rightsquigarrow \frac{1}{q} + \frac{1}{p} = 1$

Bounds on the adversarial risk:

$$R + \delta^2 N_q + \sigma^2 \leq R^{\text{adv}} \leq \left(\sqrt{R} + \delta \sqrt{N_q} \right)^2 + \sigma^2.$$

R^{adv} \rightsquigarrow Adversarial risk

R \rightsquigarrow Risk

$$N_q = E\{\|\hat{\beta}\|_q^2\}$$

\rightarrow For an ℓ_p attack $\rightsquigarrow \frac{1}{q} + \frac{1}{p} = 1$

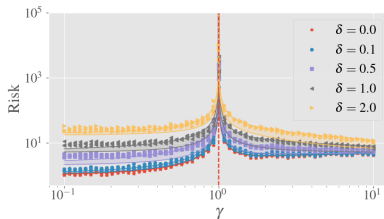
δ \rightsquigarrow Adversarial disturbance magnitude

Results

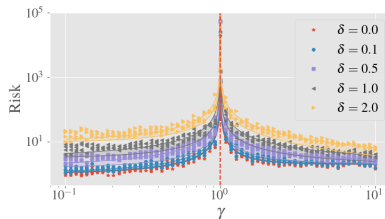
For the ℓ_2 -norm we know the asymptotic for all of the above quantities!

Results

For the ℓ_2 -norm we know the asymptotic for all of the above quantities!



(a) Uncorrelated feature

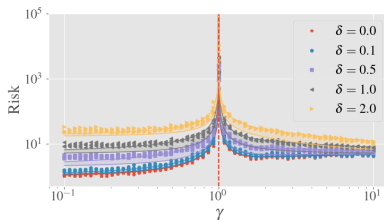


(b) Equicorrelated features ($\rho = 0.5$)

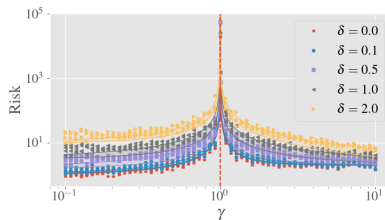
Figure: Adversarial ℓ_2 risk. The solid line gives the asymptotic risk lower and upper bound computed. The points correspond to the empirical adversarial risk for experiments.

Results

For the ℓ_2 -norm we know the asymptotic for all of the above quantities!



(a) Uncorrelated feature



(b) Equicorrelated features ($\rho = 0.5$)

Figure: Adversarial ℓ_2 risk. The solid line gives the asymptotic risk lower and upper bound computed. The points correspond to the empirical adversarial risk for experiments.

T. Hastie, A. Montanari, S. Rosset, and R. J. Tibshirani, "Surprises in High-Dimensional Ridgeless Least Squares Interpolation," arXiv:1903.08560, Nov. 2019.

- ▶ This is not in line with what is observed for l_∞ adversarial attacks ...

- ▶ This is not in line with what is observed for l_∞ adversarial attacks ...
- ▶ There adversarial performance degrades as we increase the number of features. See, for instance:

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


- ▶ This is not in line with what is observed for l_∞ adversarial attacks ...
- ▶ There adversarial performance degrades as we increase the number of features. See, for instance:
D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Ma. Robustness May Be At Odds with Accuracy. ICLR, 2019.
- ▶ Different l_p adversarial attacks may behave *qualitatively* different as we increase the number of parameters...


Thank you!

Contact info:

 **antonio.horta.ribeiro@it.uu.se**
thomas.schon@it.uu.se

 @ahortaribeiro

 antonior92.github.io
user.it.uu.se/~thosc112

 github.com/antonior92
github.com/thomasschon