On the robustness of overparametrized models

Antônio H. Ribeiro

"To think is to forget a difference, to generalize, to abstract. In the overly replete world of Funes, there were nothing but details." Funes, the Memorious (1942) Jorge Luis Borges



 μ -seminar November 12th, 2021

Model size in neural networks



Figure: Models number of parameters

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

M. Tan and Q. V. Le (2019) "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," Proceedings of the 36th International Conference on Machine Learning (ICML). PMLR, vol. 97.

Overparametrized models

The unreasonable effectiveness of overparameterized machine learning models (2021, Fall) — PhD level course (together with Dave and Per).

All material available in: https://github.com/uu-sml/seminars-overparam-ml

Inaugural paper: M. Belkin, D. Hsu, S. Ma, and S. Mandal (2019), "Reconciling modern machine-learning practice and the classical bias-variance trade-off," Proc Natl Acad Sci USA, vol. 116, no. 32, pp. 15849–15854, doi: 10.1073/pnas.1903070116.



2 / 15

Double-descent



Figure: Nonlinear ARX performance in Couple Eletric Drives benchmark.

Antônio H. Ribeiro, Johannes N. Hendriks, Adrian G. Wills, Thomas B. Schön. "Beyond Occam's Razor in System Identification: Double-Descent when Modeling Dynamics". Proceedings of the 19th IFAC Symposium on System Identification (SYSID) - IFAC-PapersOnLine (2021)

Double-descent in linear models

Estimated parameter: using train dataset (x_i, y_i) , $i = 1, \dots, n$:

Underparametrized:

$$\hat{\beta} = \arg\min_{\beta} \sum_{i} (\mathbf{y}_{i} - \mathbf{x}_{i}^{\mathsf{T}} \beta)^{2}$$

• Overparametrized:

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

subject to $y_{i} = \mathbf{x}_{i}^{\mathsf{T}}\beta$
for every i

Random features: Belking et.al. (2019) generates the features through the nonlinear mapping: $\phi : u_i \mapsto x_i$ obtained from Random Fourier Features.

Benign overfitting

- When the solution in the overparametrized region tends to the optimal one?
- Depend on the eigenvalues of the covariance matrix
 Σ = Cov [x].

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, doi: 10.1073/pnas.1907378117.

Example: Latent space model

Data is generated from a lower dimensional subspace: $\dim(z_i) < \dim(x_i)$.

$$\begin{aligned} \mathbf{x}_i &= \mathbf{W} \mathbf{z}_i + \xi_i, \\ \mathbf{y}_i &= \mathbf{\theta}^\mathsf{T} \mathbf{z}_i + \epsilon_i, \end{aligned}$$

 $W \dashrightarrow$ Projects z into a higher dimensional space $\theta \dashrightarrow$ Output is a linear combination of z, through θ $\xi \dashrightarrow$ error in variables $\epsilon \dashrightarrow$ additive noise in the output

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



Figure: Model performance vs ratio between features and datapoints

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

What when there is a distributional shift? *are they robust?*

A.H. Ribeiro, T.B. Schön, "Overparametrized Regression Under L2 Adversarial Attacks". Workshop on the Theory of Overparameterized Machine Learning (TOPML) (2021)
 A.H. Ribeiro, T.B. Schön, "On the adversarial robustness of overparametrized linear regression," Work in progress.

Adversarial Attacks



Figure: Illustration of adversarial attack.

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

Linear regression under adversarial attacks

Given a data point not seen during training (x, y).

Standard risk:

$$R = E\left\{ (\mathbf{y} - \mathbf{x}^{\mathsf{T}}\hat{\boldsymbol{\beta}})^2 | \underbrace{\mathbf{x}_i, i = 1, \cdots, n}_{\text{training inputs}} \right\}$$

Adversarial risk:

$$R_{p}^{\mathsf{adv}} = E\left\{\max_{\|\Delta x\|_{p} \leq \delta} \left(y - (x + \Delta x)^{\mathsf{T}} \hat{\beta}\right)^{2} | \underbrace{x_{i}, i = 1, \cdots, n}_{\text{training inputs}}\right\}$$

 $\Delta x \rightsquigarrow$ Adversarially generated disturbance

Can the double-descent be observed in adversarial scenarios?

1. Enlarging the function classes we are able to find models that are smoother.

M. Belkin, D. Hsu, S. Ma, and S. Mandal. Reconciling modern machine-learning practice and the classical bias-variance trade-off. Proceedings of the National Academy of Sciences, 116(32):15849–15854, Aug. 2019. ISSN 0027-8424, 1091-6490. doi:10.1073/pnas.1903070116

S. Bubeck and M. Sellke. A Universal Law of Robustness via Isoperimetry. arXiv:2105.12806, June2021

2. Robustness-accuracy tradeoff

A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards Deep Learning Models Resistantto Adversarial Attacks.Proceedings of the International Conference for Learning Representations (ICLR), 2018.

Adversarial attacks

Geometry of the adversarial attack for dim(x) = 2:



Two types of behavior:



 μ -Seminar

Methods

Bounds on the adversarial risk:

$$R + \delta^2 L_q \leqslant R_p^{\text{adv}} \leqslant \left(\sqrt{R} + \delta\sqrt{L_q}\right)^2.$$

 $\begin{array}{l} R^{\text{adv}} & \leadsto & \text{Adversarial risk} \\ R & \leadsto & \text{Risk} \\ L_q &= E\{\|\hat{\beta}\|_q^2\} \\ & \rightarrow \text{ For an } \ell_p \text{ attack } \leadsto & \frac{1}{q} + \frac{1}{p} = 1 \\ \delta & \leadsto & \text{Adversarial disturbance magnitude} \end{array}$

Results (ℓ_2 attacks)

For the $\ell_2\text{-norm}$ 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.

Other ℓ_p attacks

$$\begin{cases} L_2 \leq L_q \leq (m)^{1-2/p}L_2, & 2
$$m \xrightarrow{\text{momentary}} \text{number of features}$$

$$L_q = E\{\|\hat{\beta}\|_q^2\}$$
For an ℓ_p attack $\xrightarrow{\text{momentary}} \frac{1}{q} + \frac{1}{p} = 1$
Figure: Adversarial risk with$$

Figure: Adversarial risk with bounds.

101

Conclusion and next steps

- Our framework gives flexibility to experiment diverse scenarios and reuse previous literature.
- We are now focusing on the question when the upper bound in the experiments above is tight.
- It depends on the geometry of the eigenvectors of Σ.



Figure: Example with isotropic features