# Mathematics of Data: From Theory to Computation

Prof. Volkan Cevher
*volkan.cevher@epfl.ch*

*Lecture 7: A mathematical introduction to Deep Learning*

Laboratory for Information and Inference Systems (LIONS)
École Polytechnique Fédérale de Lausanne (EPFL)

**EE-556** (Fall 2021)

# License Information for Mathematics of Data Slides

- This work is released under a [Creative Commons License](#) with the following terms:
- **Attribution**
  - The licensor permits others to copy, distribute, display, and perform the work. In return, licensees must give the original authors credit.
- **Non-Commercial**
  - The licensor permits others to copy, distribute, display, and perform the work. In return, licensees may not use the work for commercial purposes – unless they get the licensor's permission.
- **Share Alike**
  - The licensor permits others to distribute derivative works only under a license identical to the one that governs the licensor's work.
- [Full Text of the License](#)

**Outline**

- This class
  - ▶ Introduction to Deep Learning
  - ▶ The Deep Learning Paradigm
  - ▶ Challenges in Deep Learning Theory and Applications
  - ▶ Introduction to Generalization error bounds
    - ▶ Uniform Convergence and Rademacher Complexity
  - ▶ Generalization in Deep Learning (Part 1)
- Next class
  - ▶ Generalization in Deep Learning (Part 2)

## Remark about notation

The Deep Learning literature might use a different notation:

|  | Our lectures | DL literature |
|---|:---:|:---:|
| data/sample | $\mathbf{a}$ | $\mathbf{x}$ |
| label | $b$ | $y$ |
| bias | $\mu$ | $b$ |
| weight | $\mathbf{x}, \mathbf{X}$ | $\mathbf{w}, \mathbf{W}$ |

## Power of linear classifiers–I

### Problem (Recall: Logistic regression)

*Given a sample vector $\mathbf{a}_i \in \mathbb{R}^d$ and a binary class label $b_i \in \{-1, +1\}$ $(i = 1, \ldots, n)$, we define the conditional probability of $b_i$ given $\mathbf{a}_i$ as follows:*

$$\mathbb{P}(b_i | \mathbf{a}_i, \mathbf{x}) \propto 1/(1 + e^{-b_i \langle \mathbf{x}, \mathbf{a}_i \rangle}),$$

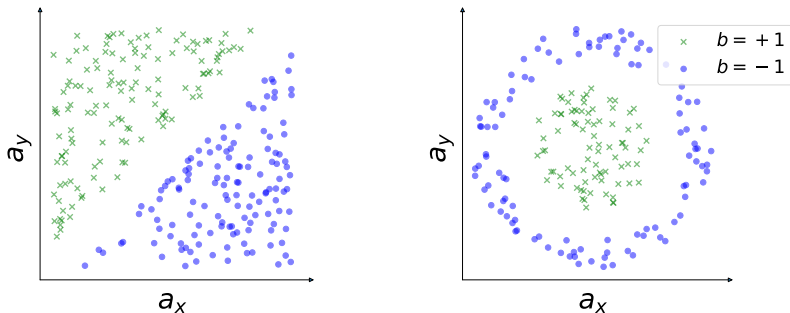*where $\mathbf{x} \in \mathbb{R}^d$ is some weight vector.*



Figure: Linearly separable versus nonlinearly separable dataset

# Power of linear classifiers–II

- Lifting dimensions to the rescue
  - ▶ Convex optimization objective
  - ▶ Side effect: The curse-of-dimensionality
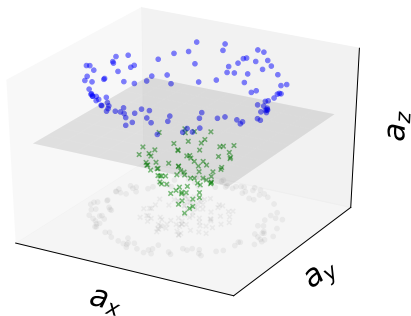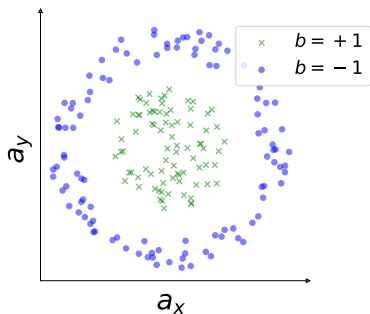  - ▶ Possible to avoid via kernel methods, such as SVMs



Figure: Non-linearly separable data (left). Linearly separable in $\mathbb{R}^3$ via $\mathbf{a}_z = \sqrt{\mathbf{a}_x^2 + \mathbf{a}_y^2}$ (right).

## An important alternative for non-linearly separable data

**1-hidden-layer neural network with $m$ neurons (fully-connected architecture):**

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m, \mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \to \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) :=$$

input
↓

$$\begin{bmatrix} \\ \mathbf{a} \\ \\ \end{bmatrix}$$

# An important alternative for non-linearly separable data

**1-hidden-layer neural network with $m$ neurons (fully-connected architecture):**

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m, \mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \to \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \begin{bmatrix} \overset{\text{weight}}{\downarrow} \\ \mathbf{X}_1 \end{bmatrix} \begin{bmatrix} \overset{\text{input}}{\downarrow} \\ \mathbf{a} \end{bmatrix}$$

# An important alternative for non-linearly separable data

**1-hidden-layer neural network with $m$ neurons (fully-connected architecture):**

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m, \mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \to \mathbb{R}$

$$
h_{\mathbf{x}}(\mathbf{a}) := \left[ \begin{array}{c} \overset{\text{weight}}{\underset{\downarrow}{}} \\ \mathbf{X}_1 \end{array} \right] \left[ \overset{\text{input}}{\underset{\downarrow}{\mathbf{a}}} \right] + \left[ \overset{\text{bias}}{\underset{\downarrow}{\mu_1}} \right]
$$

# An important alternative for non-linearly separable data

**1-hidden-layer neural network with $m$ neurons (fully-connected architecture):**

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m, \mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \to \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \underset{\substack{\uparrow \\ \text{activation}}}{\sigma} \left( \left[ \quad \underset{\substack{\uparrow \\ \text{weight}}}{\mathbf{X}_1} \quad \right] \left[ \underset{\substack{\uparrow \\ \text{input}}}{\mathbf{a}} \right] + \left[ \underset{\substack{\uparrow \\ \text{bias}}}{\mu_1} \right] \right)$$

hidden layer = learned features

**An important alternative for non-linearly separable data**

**1-hidden-layer neural network with $m$ neurons (fully-connected architecture):**

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m, \mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \to \mathbb{R}$



$$h_{\mathbf{x}}(\mathbf{a}) := \left[\;\; \mathbf{X}_2 \;\;\right] \underset{\uparrow}{\overset{\text{activation}}{\sigma}} \left(\left[\;\; \mathbf{X}_1 \;\;\right]\underset{\uparrow}{\overset{\text{weight}}{}}\left[\mathbf{a}\right]\underset{\uparrow}{\overset{\text{input}}{}} + \left[\mu_1\right]\underset{\uparrow}{\overset{\text{bias}}{}}\right)$$

hidden layer = learned features

# An important alternative for non-linearly separable data

**1-hidden-layer neural network with $m$ neurons (fully-connected architecture):**

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m, \mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \to \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \begin{bmatrix} \mathbf{X}_2 \end{bmatrix} \underset{\uparrow}{\sigma} \left( \begin{bmatrix} \mathbf{X}_1 \end{bmatrix} \underset{\uparrow}{\begin{bmatrix} \mathbf{a} \end{bmatrix}} + \underset{\uparrow}{\begin{bmatrix} \mu_1 \end{bmatrix}} \right) + \underset{\uparrow}{\begin{bmatrix} \mu_2 \end{bmatrix}}$$

activation · weight · input · bias · bias

$\underbrace{\qquad\qquad\qquad\qquad}_{\text{hidden layer} \,=\, \text{learned features}}$

## An important alternative for non-linearly separable data

**1-hidden-layer neural network with $m$ neurons (fully-connected architecture):**

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m, \mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \to \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \left[ \begin{array}{c} \mathbf{X}_2 \end{array} \right] \underset{\underset{\textcolor{red}{\text{activation}}}{\downarrow}}{\textcolor{red}{\sigma}} \left( \left[ \begin{array}{c} \mathbf{X}_1 \end{array} \right] \underset{\underset{\text{weight}}{\downarrow}}{} \left[ \textcolor{green}{\mathbf{a}} \right] \underset{\underset{\text{input}}{\downarrow}}{} + \left[ \textcolor{blue}{\mu_1} \right] \underset{\underset{\textcolor{blue}{\text{bias}}}{\downarrow}}{} \right) + \left[ \textcolor{blue}{\mu_2} \right] \underset{\underset{\textcolor{blue}{\text{bias}}}{\downarrow}}{}, \qquad \mathbf{x} := [\mathbf{X}_1, \mathbf{X}_2, \mu_1, \mu_2]$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}_{\textcolor{red}{\text{hidden layer} = \text{learned features}}}$$

## An important alternative for non-linearly separable data

**1-hidden-layer neural network with $m$ neurons (fully-connected architecture):**

- Parameters: $\mathbf{X}_1 \in \mathbb{R}^{m \times d}$, $\mathbf{X}_2 \in \mathbb{R}^{c \times m}$ (weights), $\mu_1 \in \mathbb{R}^m$, $\mu_2 \in \mathbb{R}^c$ (biases)
- Activation function: $\sigma : \mathbb{R} \to \mathbb{R}$

$$h_{\mathbf{x}}(\mathbf{a}) := \left[ \quad \mathbf{X}_2 \quad \right] \underset{\sigma}{\overset{\text{activation}}{\downarrow}} \left( \left[ \quad \mathbf{X}_1 \quad \right] \underset{[\mathbf{a}]}{\overset{\text{input}}{\downarrow}} + \underset{[\mu_1]}{\overset{\text{bias}}{\downarrow}} \right) + \underset{[\mu_2]}{\overset{\text{bias}}{\downarrow}}, \qquad \mathbf{x} := [\mathbf{X}_1, \mathbf{X}_2, \mu_1, \mu_2]$$

(with $\mathbf{X}_1$ labeled **weight**)

underbrace: hidden layer = learned features

recursively repeat activation + affine transformation to obtain "deeper" networks.

# Why neural networks?: An approximation theoretic motivation

## Theorem (Universal approximation [4])

*Let $\sigma(\cdot)$ be a nonconstant, bounded, and increasing continuous function. Let $I_d = [0,1]^d$. The space of continuous functions on $I_d$ is denoted by $\mathcal{C}(I_d)$.*

*Given $\epsilon > 0$ and $g \in \mathcal{C}(I_d)$ there exists a 1-hidden-layer network $h$ with $m$ neurons such that $h$ is an $\epsilon$-approximation of $g$, i.e.,*

$$\sup_{\mathbf{a} \in I_d} |g(\mathbf{a}) - h(\mathbf{a})| \le \epsilon$$

## Caveat

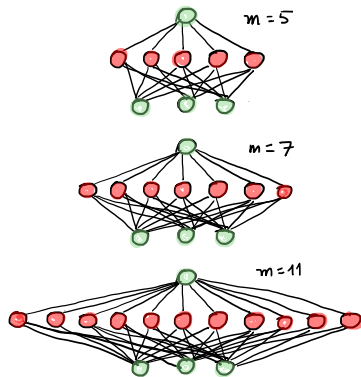The number of neurons $m$ needed to approximate some function $g$ can be arbitrarily large!



Figure: networks of increasing width

# Why were NNs not popular before 2010?

- too big to optimize!
- did not have enough data
- could not find the optimum via algorithms

# Why were NNs not popular before 2010?

- ▶ too big to optimize!
- ▶ did not have enough data
- ▶ could not find the optimum via algorithms

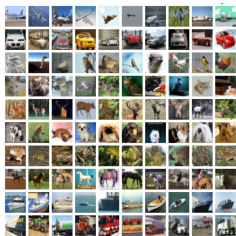# Supervised learning: Multi-class classification



Figure: CIFAR10 dataset: 60000 32x32 color images (3 channels) from 10 classes



Figure: Imagenet dataset: 14 million color images (varying resolution, 3 channels) from 21K classes

## Goal

Image-label pairs $(\mathbf{a}, b) \subseteq \mathbb{R}^d \times \{1, \ldots, c\}$ follow an unknown distibution $\mathbb{P}$. Find $h : \mathbb{R}^d \to \{1, \ldots, c\}$ with minimum *misclassification probability*

$$\min_{h \in \mathcal{H}} \mathbb{P}(h(\mathbf{a}) \neq b)$$

## 2010-today: Deep Learning becomes popular again



Figure: Error rate on the ImageNet challenge, for different network architectures.

# 2010-today: Deep Learning becomes popular again



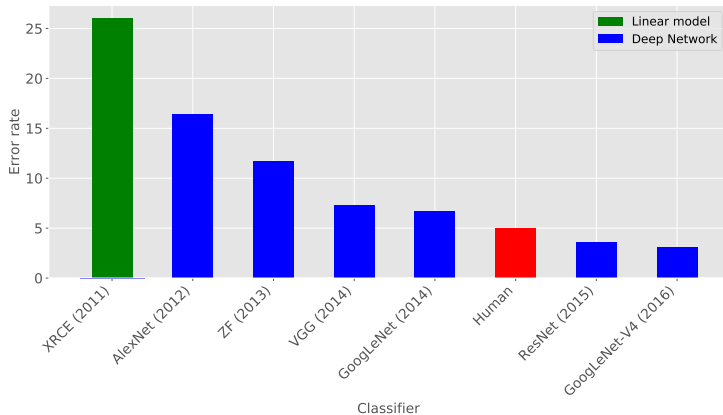Figure: Error rate on the ImageNet challenge, for different network architectures.[?, ?]

# Convolutional architectures in Computer Vision tasks



Figure: "Locality" Structure of a 2D deep convolutional neural network.

# Inductive Bias: Why convolution works so well in Computer Vision tasks?



| | |
|---|---|
| $h^\circ$ | true unknown function |
| $\mathcal{H}$ | space of all functions |
| $\mathcal{H}_{\text{fc}}^p$ | fully-connected networks with $p$ parameters |
| $\mathcal{H}_{\text{conv}}^p$ | convolutional networks with $p$ parameters |

# 2010-today: Size of neural networks grows exponentially!



Figure: Number of parameters in Language models based on Deep Learning.

# The era of model scaling



Figure: Scale of most recent SotA models across modalities.
Data from https://lair.lighton.ai/akronomicon/:

# The Landscape of ERM with multilayer networks

## Recall: Empirical risk minimization (ERM)

Let $h_{\mathbf{x}} : \mathbb{R}^n \to \mathbb{R}$ be network and let $\{(\mathbf{a}_i, b_i)\}_{i=1}^n$ be a sample with $b_i \in \{-1, 1\}$ and $\mathbf{a}_i \in \mathbb{R}^n$. The *empirical risk minimization* (ERM) is defined as follows

$$\min_{\mathbf{x}} \left\{ R_n(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n L(h_{\mathbf{x}}(\mathbf{a}_i), b_i) \right\} \tag{1}$$

where $L(h_{\mathbf{x}}(\mathbf{a}_i), b_i)$ is the loss on the sample $(\mathbf{a}_i, b_i)$ and $\mathbf{x}$ are the parameters of the network.

## Some frequently used loss functions

- $L(h_{\mathbf{x}}(\mathbf{a}), b) = \log(1 + \exp(-b \cdot h_{\mathbf{x}}(\mathbf{a})))$ (logistic loss)
- $L(h_{\mathbf{x}}(\mathbf{a}), b) = (b - h_{\mathbf{x}}(\mathbf{a}))^2$ (squared error)
- $L(h_{\mathbf{x}}(\mathbf{a}), b) = \max(0, 1 - b \cdot h_{\mathbf{x}}(\mathbf{a}))$ (hinge loss)

# The Landscape of ERM with multilayer networks



Figure: convex (left) vs non-convex (right) optimization landscape

Conventional wisdom in ML until 2010:
Simple models + simple errors

# The deep learning paradigm



(a) Massive datasets

(b) Inductive bias from large and complex architectures

(c) ERM using stochastic non-convex first-order optimization algorithms (SGD)

Figure: Most common components in a Deep Learning Pipeline

# Challenges in DL/ML applications: Robustness (I)



(a) Turtle classified as rifle [1].



(b) Stop sign classified as 45 mph sign [6].

Figure: Natural or human-crafted modifications that trick neural networks used in computer vision tasks

# Challenges in DL/ML applications: Robustness (II)



(a) Linear classifier on data distributed on a sphere

$$h_{\mathbf{x}}(\mathbf{a} + \epsilon) \approx h_{\mathbf{x}}(\mathbf{a}) + \langle \epsilon, \nabla h_{\mathbf{x}}(\mathbf{a}) \rangle$$
$$|h_{\mathbf{x}}(\mathbf{a} + \epsilon) - h_{\mathbf{x}}(\mathbf{a})| \leq \|\epsilon\| \|\nabla h_{\mathbf{x}}(\mathbf{a})\|$$

(b) Concentration of measure phenomenon on high dimensions

Figure: Understanding the robustness of a classifier in high-dimensional spaces [12]

# Challenges in DL/ML applications: Robustness (References I)

1. Madry, Aleksander and Makelov, Aleksandar and Schmidt, Ludwig and Tsipras, Dimitris and Vladu, Adrian. *Towards Deep Learning Models Resistant to Adversarial Attacks*. ICLR 2018.

2. Raghunathan, A., Steinhardt, J., and Liang, P. S. *Semidefinite relaxations for certifying robustness to adversarial examples*. Neurips 2018.

3. Wong, E. and Kolter, Z. (2018). *Provable defenses against adversarial examples via the convex outer adversarial polytope*. ICML 2018.

4. Huang, X., Kwiatkowska, M., Wang, S., and Wu, M. *Safety verification of deep neural networks*. Computer Aided Verification 2017.

5. Athalye, A., et al. *Synthesizing robust adversarial examples*. International conference on machine learning. PMLR, 2018.

6. Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., and Song, D. *Robust physical-world attacks on deep learning visual classification*. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 1625-1634). 2018.

7. Shafahi A., Ronny Huang, W., Studer, C., Feizi, S. and Goldstein, T. *Are adversarial examples inevitable?*. International Conference on Learning Representations. 2019.

## Challenges in DL/ML applications: Robustness (References II)

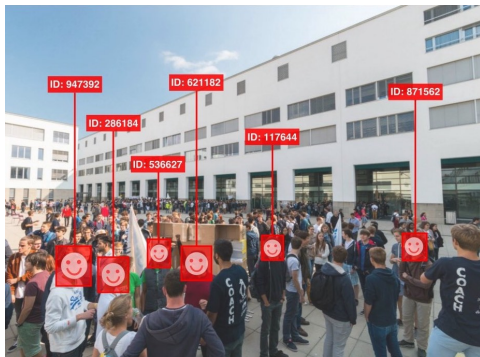1. Z. Charles, H. Rosenberg, and D. Papailiopoulos, A Geometric Perspective on the Transferability of Adversarial Directions, in Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics, Apr. 2019, pp. 1960–1968.

2. J. Gilmer, N. Ford, N. Carlini, and E. Cubuk, Adversarial Examples Are a Natural Consequence of Test Error in Noise, in Proceedings of the 36th International Conference on Machine Learning, May 2019, pp. 2280–2289.

3. A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, Adversarial Examples Are Not Bugs, They Are Features, in Advances in Neural Information Processing Systems, 2019, vol. 32.

4. A. Fawzi, H. Fawzi, and O. Fawzi, Adversarial vulnerability for any classifier, in Advances in Neural Information Processing Systems, 2018, vol. 31.

5. L. Schmidt, S. Santurkar, D. Tsipras, K. Talwar, and A. Ma, *Adversarially Robust Generalization Requires More Data*, p. 13.

6. J.-H. Jacobsen, J. Behrmann, R. Zemel, and M. Bethge, Excessive Invariance Causes Adversarial Vulnerability, presented at the International Conference on Learning Representations, Sep. 2018.

7. F. Tramer, J. Behrmann, N. Carlini, N. Papernot, and J.-H. Jacobsen, Fundamental Tradeoffs between Invariance and Sensitivity to Adversarial Perturbations, in Proceedings of the 37th International Conference on Machine Learning, Nov. 2020, pp. 9561–9571.

8. C. Xie and A. Yuille, Intriguing Properties of Adversarial Training at Scale, presented at the International Conference on Learning Representations, Sep. 2019.

# Challenges in DL/ML applications: Robustness (References III)

1. L. Chen, Y. Min, M. Zhang, and A. Karbasi, *More Data Can Expand the Generalization Gap Between Adversarially Robust and Standard Models,* p. 11.
2. F. Tramèr, N. Carlini, W. Brendel, and A. Ma, *On Adaptive Attacks to Adversarial Example Defenses,* p. 37.
3. D. Krueger et al., Out-of-Distribution Generalization via Risk Extrapolation (REx), arXiv:2003.00688 [cs, stat], Feb. 2021,
4. T. Pang, K. Xu, Y. Dong, C. Du, N. Chen, and J. Zhu, *Rethinking Softmax Cross-Entropy Loss For Adversarial Robustness,* p. 19, 2020.
5. R. Bhattacharjee, S. Jha, and K. Chaudhuri, Sample Complexity of Robust Linear Classification on Separated Data, in Proceedings of the 38th International Conference on Machine Learning, Jul. 2021, pp. 884–893.
6. R. Geirhos et al., *Shortcut Learning in Deep Neural Networks,* Nat Mach Intell, vol. 2, no. 11, pp. 665–673, Nov. 2020, doi: 10.1038/s42256-020-00257-z.
7. Y. Min, L. Chen, and A. Karbasi, The Curious Case of Adversarially Robust Models: More Data Can Help, Double Descend, or Hurt Generalization, arXiv:2002.11080 [cs, stat], Jun. 2020
8. E. Rosenfeld, P. K. Ravikumar, and A. Risteski, The Risks of Invariant Risk Minimization, presented at the International Conference on Learning Representations, Sep. 2020.
9. T. Li, A. Beirami, M. Sanjabi, and V. Smith, Tilted Empirical Risk Minimization, presented at the International Conference on Learning Representations, Sep. 2020.
10. A. D'Amour et al., Underspecification Presents Challenges for Credibility in Modern Machine Learning, arXiv:2011.03395 [cs, stat], Nov. 2020,

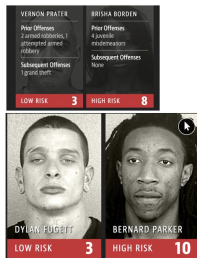# Challenges in DL/ML applications: Surveillance/Privacy/Manipulation



Figure: Political and societal concerns about some DL/ML applications

# Challenges in DL/ML applications: Surveillance/Privacy/Manipulation (References)

1. Dwork, C., and Roth, A. *The Algorithmic Foundations of Differential Privacy*. Foundations and Trends in Theoretical Computer Science, 9, 2013.
2. Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. *Deep learning with differential privacy*. In Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (pp. 308-318). 2016.
3. Sreenu, G., Saleem Durai, M.A. *Intelligent video surveillance: a review through deep learning techniques for crowd analysis*. J Big Data 6, 48. 2019.
4. O'Neil, C., *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy* Broadway Books, (2016);
5. Ali, R. E., So, J., Avestimehr, A. S. *On polynomial approximations for privacy-preserving and verifiable relu networks*. arXiv preprint arXiv:2011.05530.
6. Bagdasaryan, E., Poursaeed, O., Shmatikov, V. *Differential privacy has disparate impact on model accuracy*. In Neural Information Processing Systems (NeurIPS), 2019.
7. Bu, Z., Dong, J., Long, Q., Su, W. J. *Deep learning with Gaussian differential privacy*. Harvard data science review, 2020.
8. Pujol, D., McKenna, R., Kuppam, S., Hay, M., Machanavajjhala, A., Miklau, G. (textitFair decision making using privacy-protected data . In Proceedings of Fairness, Accountability, and Transparency (FAT), 2020.
9. Tran, C., Fioretto, F., Van Hentenryck, P. *Differentially private and fair deep learning: A lagrangian dual approach*. arXiv preprint arXiv:2009.12562.
10. Xu, D., Du, W., Wu, X. *Removing disparate impact on model accuracy in differentially private stochastic gradient descent*. In Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, 2021.

# Challenges in DL/ML applications: Fairness



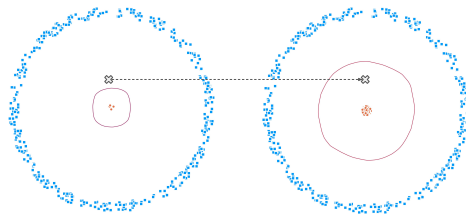(a) Racist classifier

(b) Effect of unbalanced data

Figure: Unfair classifiers due to biased or unbalanced datasets/algorithms

# Challenges in DL/ML applications: Fairness (References)

1. Barocas, S. Hardt, M. Narayanan, Arvind. *Fairness in Machine Learning Limitations and Opportunities*. `https://fairmlbook.org/pdf/fairmlbook.pdf` 2020.

2. Hardt, M. *How Big Data Is Unfair*. `https://medium.com/@mrtz/how-big-data-is-unfair-9aa544d739de` 2014.

3. Munoz, C., Smith, M., and Patil, D. *Big Data: A Report on Algorithmic Systems, Opportunity, and Civil Rights*. Executive Office of the President. The White House, 2016.

4. Campolo, A., Sanfilippo, M., Whittaker, M., Crawford, K. *AI Now 2017 Report*. AI Now Institute at New York University, 2017.

5. Friedman, B. and Nissenbaum, H. *Bias in Computer Systems*. ACM Transactions on Information Systems (TOIS) 14, no. 3. 1996: 330–47.

6. Pedreshi, D., Ruggieri, S. and Turini, F. *Discrimination-Aware Data Mining*. Proc. 14th SIGKDD. ACM 2008.

7. Noble, S.U. *Algorithms of Oppression: How Search Engines Reinforce Racism*. NYU Press. 2018.

8. Rolf, E., Simchowitz, M., Dean, S., Liu, L. T., Bjorkegren, D., Hardt, M., Blumenstock, J. *Balancing competing objectives with noisy data: Score-based classifiers for welfare-aware machine learning*. In International Conference on Machine Learning (ICML), 2021.

9. Hanna, A., Denton, E., Smart, A., Smith-Loud, J. *Towards a critical race methodology in algorithmic fairness*. In Proceedings of Fairness, Accountability, and Transparency (FAT), 2020.

10. Wang, A., Russakovsky, O. *Directional bias amplification*. In International Conference on Machine Learning (ICML), 2021.

# Challenges in DL/ML applications: Fairness (References)

1. Yang, K., Qinami, K., Fei-Fei, L., Deng, J., Russakovsky, O. *Towards fairer datasets: Filtering and balancing the distribution of the people subtree in the imagenet hierarchy*. In Proceedings of Fairness, Accountability, and Transparency (FAT), 2020.

2. Mitchell, M., Baker, D., Moorosi, N., Denton, E., Hutchinson, B., Hanna, A., Morgenstern, J. *Diversity and inclusion metrics in subset selection*. In Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society, 2020.

3. Obermeyer, Z., Powers, B., Vogeli, C., Mullainathan, S. *Dissecting racial bias in an algorithm used to manage the health of populations*. Science, 366(6464), 447-453, 2019.

4. Jia, S., Meng, T., Zhao, J., Chang, K. W. *Mitigating gender bias amplification in distribution by posterior regularization*. Annual Meeting of the Association for Computational Linguistics (ACL), 2020.

5. Jalal, A., Karmalkar, S., Hoffmann, J., Dimakis, A., Price, E. *Fairness for Image Generation with Uncertain Sensitive Attributes*. In International Conference on Machine Learning (ICML), 2021.

6. Hoyle, A., Wallach, H., Augenstein, I., Cotterell, R. *Unsupervised discovery of gendered language through latent-variable modeling*. arXiv preprint arXiv:1906.04760.

7. Ramaswamy, V. V., Kim, S. S., Russakovsky, O. *Fair attribute classification through latent space de-biasing*. In Proceedings of the Conference on Computer Vision and Pattern Recognition, 2021.

8. Jacobs, A. Z., Wallach, H. *Measurement and fairness*. In Proceedings of Fairness, Accountability, and Transparency (FAT), 2021.

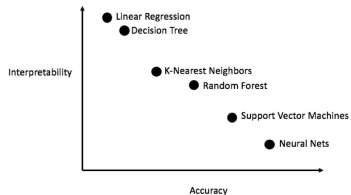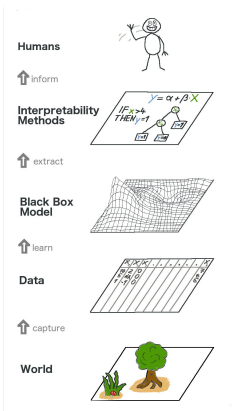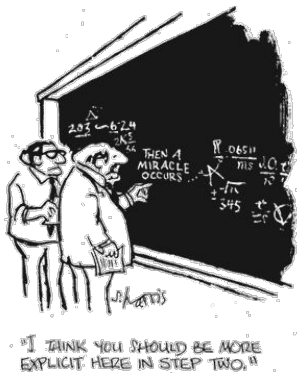# Challenges in DL/ML applications: Interpretability



Figure: Performance vs Interpretability trade-offs in DL/ML

# Challenges in DL/ML applications: Interpretability (References I)

1. Baehrens, David and Schroeter, Timon and Harmeling, Stefan and Kawanabe, Motoaki and Hansen, Katja and Mueller, Klaus-Robert. Simonyan, Karen and Vedaldi, Andrea and Zisserman, Andrew. How to Explain Individual Classification Decisions. JMLR 2010.

2. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. arXiv e-prints. arXiv:1312.6034. 2013.

3. Ribeiro, Marco and Singh, Sameer and Guestrin, Carlos. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. KDD 2016.

4. Sundararajan, Mukund and Taly, Ankur and Yan, Qiqi. Axiomatic Attribution for Deep Networks. ICML 2017.

5. Shrikumar, Avanti and Greenside, Peyton and Kundaje, Anshul. Learning Important Features Through Propagating Activation Differences. ICML 2017.

6. C. Rudin, C. Chen, Z. Chen, H. Huang, L. Semenova, and C. Zhong, Interpretable Machine Learning: Fundamental Principles and 10 Grand Challenges arXiv:2103.11251 [cs, stat], Jul. 2021

7. L. Semenova, C. Rudin, and R. Parr, A study in Rashomon curves and volumes: A new perspective on generalization and model simplicity in machine learning arXiv:1908.01755 [cs, stat], Apr. 2021
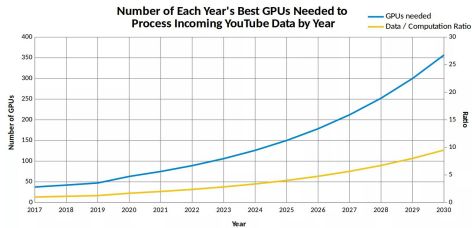
# Challenges in DL/ML applications: Interpretability (References II)

1. S. Hooker, D. Erhan, P.-J. Kindermans, and B. Kim, *A Benchmark for Interpretability Methods in Deep Neural Networks* p. 12.

2. F. Poursabzi-Sangdeh, D. G. Goldstein, J. M. Hofman, J. W. Wortman Vaughan, and H. Wallach, Manipulating and Measuring Model Interpretability in Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems, New York, NY, USA: Association for Computing Machinery, 2021, pp. 1–52.

3. H. Kaur, H. Nori, S. Jenkins, R. Caruana, H. Wallach, and J. Wortman Vaughan, *Interpreting Interpretability: Understanding Data Scientists' Use of Interpretability Tools for Machine Learning* in Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems, New York, NY, USA, Apr. 2020, pp. 1–14. doi: 10.1145/3313831.3376219.

4. P. Hase and M. Bansal, *Evaluating Explainable AI: Which Algorithmic Explanations Help Users Predict Model Behavior?* in Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Online, 2020, pp. 5540–5552. doi: 10.18653/v1/2020.acl-main.491.

5. C. Rudin and J. Radin, *Why Are We Using Black Box Models in AI When We Don't Need To? A Lesson From An Explainable AI Competition* Harvard Data Science Review, vol. 1, no. 2, Nov. 2019, doi: 10.1162/99608f92.5a8a3a3d.

6. C. Rudin, *Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead* Nat Mach Intell, vol. 1, no. 5, pp. 206–215, May 2019, doi: 10.1038/s42256-019-0048-x.

7. S. Barocas, M. Hardt, and A. Narayanan, *Fairness in machine learning* Nips tutorial, vol. 1, p. 2017, 2017.

8. J. Lee, S. Park, and J. Shin, Learning Bounds for Risk-sensitive Learning arXiv:2006.08138 [cs, stat], Jan. 2021,

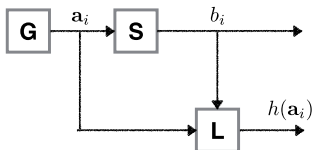# Challenges in DL/ML applications: Energy efficiency and cost



(a)

(b)

Figure: Efficiency and Scalability concerns in DL/ML

# Challenges in DL/ML applications: Energy efficiency and cost (References)

1. García-Marín, E., Rodrigues, C. F., Riley, G., and Grahn, H. *Estimation of energy consumption in machine learning*. Journal of Parallel and Distributed Computing, 134, 75-88. 2019.

2. Strubell, E., Ganesh, A., and McCallum, A. *Energy and policy considerations for deep learning in NLP*. arXiv preprint arXiv:1906.02243. 2019.

3. Goel, A., Tung, C., Lu, Y. H., and Thiruvathukal, G. K. *A Survey of Methods for Low-Power Deep Learning and Computer Vision*. arXiv preprint arXiv:2003.11066. 2020.

4. Conti, F., Rusci, M., and Benini, L. *The Memory Challenge in Ultra-Low Power Deep Learning*. In NANO-CHIPS 2030 (pp. 323-349). Springer, Cham. 2020.

5. J. Launay, I. Poli, F. Boniface, and F. Krzakala, Direct Feedback Alignment Scales to Modern Deep Learning Tasks and Architectures NeurIPS 2020

6. J. Launay et al., Hardware Beyond Backpropagation: a Photonic Co-Processor for Direct Feedback Alignment arXiv:2012.06373

7. E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? in Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, New York, NY, USA, Mar. 2021, pp. 610–623. doi: 10.1145/3442188.3445922.

# Peeling the onion



## Models

Let $d(\cdot,\cdot): \mathcal{H}^\circ \times \mathcal{H}^\circ \to \mathbb{R}^+$ be a metric in an extended function space $\mathcal{H}^\circ$ that includes $\mathcal{H}$; i.e., $\mathcal{H} \subseteq \mathcal{H}^\circ$. Let

1. $h^\circ \in \mathcal{H}^\circ$ be the true, expected risk minimizing model
2. $h^\natural \in \mathcal{H}$ be the solution under the assumed function class $\mathcal{H} \subseteq \mathcal{H}^\circ$
3. $h^\star \in \mathcal{H}$ be the estimator solution
4. $h^t \in \mathcal{H}$ be the numerical approximation of the algorithm at time $t$

## Practical performance

$$\underbrace{d(h^t, h^\circ)}_{\bar{\varepsilon}(t,n)} \leq \underbrace{d(h^t, h^\star)}_{\text{optimization error}} + \underbrace{d(h^\star, h^\natural)}_{\text{statistical error}} + \underbrace{d(h^\natural, h^\circ)}_{\text{model error}},$$

where $\bar{\varepsilon}(t,n)$ denotes the total error of the Learning Machine. We can try to

1. reduce the optimization error with computation
2. reduce the statistical error with more data samples, with better estimators, and with prior information
3. reduce the model error with flexible or universal representations

# Estimation of parameters vs estimation of risk

## Recall the general setting

Let $R(h_{\mathbf{x}}) = \mathbb{E}L(h_{\mathbf{x}}(\mathbf{a}), b)$ be the risk function and
$R_n(h_{\mathbf{x}}) = \frac{1}{n} \sum_{i=1}^{n} L(h_{\mathbf{x}}(\mathbf{a}_i), b_i)$ be the empirical estimate.
Let $\mathcal{X} \subseteq \mathcal{X}^{\circ}$ be parameter domains, where $\mathcal{X}$ is known. Define

1. $\mathbf{x}^{\circ} \in \arg\min_{\mathbf{x} \in \mathcal{X}^{\circ}} R(h_{\mathbf{x}})$: true minimum risk model
2. $\mathbf{x}^{\natural} \in \arg\min_{\mathbf{x} \in \mathcal{X}} R(h_{\mathbf{x}})$: assumed minimum risk model
3. $\mathbf{x}^{\star} \in \arg\min_{\mathbf{x} \in \mathcal{X}} R_n(h_{\mathbf{x}})$: ERM solution
4. $\mathbf{x}^t$: numerical approximation of $\mathbf{x}^{\star}$ at time $t$

## Nomenclature

| | |
|---:|:---|
| $R_n(\cdot)$ | training error |
| $R(\cdot)$ | test error |
| $R(\mathbf{x}^{\natural}) - R(\mathbf{x}^{\circ})$ | modeling error |
| $R(\mathbf{x}^{\star}) - R(\mathbf{x}^{\natural})$ | excess risk |
| $\sup_{\mathbf{x} \in \mathcal{X}} |R(\mathbf{x}) - R_n(\mathbf{x})|$ | generalization error |
| $R_n(\mathbf{x}^t) - R_n(\mathbf{x}^{\star})$ | optimization error |

| | $\mathcal{X} \to \mathcal{X}^{\circ}$ | $n \uparrow$ | $p \uparrow$ |
|---|:---:|:---:|:---:|
| Training error | ↘ | ↗ | ↘ |
| Excess risk | ↗ | ↘ | ↗ |
| Generalization error | ↗ | ↘ | ↗ |
| Modeling error | ↘ | = | ↔ |
| Time | ↗ | ↗ | ↗ |

## What theoretical challenges in Deep Learning will we study?



**Models**

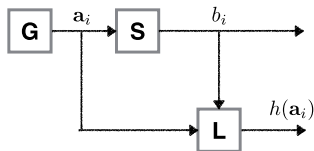Let $\mathcal{X} \subseteq \mathcal{X}^\circ$ be parameter domains, where $\mathcal{X}$ is known. Define

1. $\mathbf{x}^\circ \in \arg\min_{\mathbf{x} \in \mathcal{X}^\circ} R(h_\mathbf{x})$: true minimum risk model
2. $\mathbf{x}^\natural \in \arg\min_{\mathbf{x} \in \mathcal{X}} R(h_\mathbf{x})$: assumed minimum risk model
3. $\mathbf{x}^\star \in \arg\min_{\mathbf{x} \in \mathcal{X}} R_n(h_\mathbf{x})$: ERM solution
4. $\mathbf{x}^t$: numerical approximation of $\mathbf{x}^\star$ at time $t$

**Practical performance in Deep Learning**

$$\underbrace{R(\mathbf{x}^t) - R(\mathbf{x}^\circ)}_{\bar{\varepsilon}(t,n)} \leq \underbrace{R_n(\mathbf{x}^t) - R_n(\mathbf{x}^\star)}_{\text{optimization error}} + \underbrace{2 \sup_{\mathbf{x} \in \mathcal{X}} |R(\mathbf{x}) - R_n(\mathbf{x})|}_{\text{generalization error}} + \underbrace{R(\mathbf{x}^\natural) - R(\mathbf{x}^\circ)}_{\text{model error}}$$

where $\bar{\varepsilon}(t,n)$ denotes the total error of the Learning Machine. In Deep Learning applications

1. Optimization error is almost zero, in spite of non-convexity. $\Rightarrow$ lecture 9
2. Generalization error is usually small, but theory is lacking. $\Rightarrow$ lecture 7 (this one) and lecture 8
3. Large architectures + inductive bias might lead to small model error.

# Generalization error bounds

Goal: Obtain generalization bounds for multi-layer, fully-connected neural networks

○ We want to find high-probability upper bounds for the quantity

$$\sup_{\mathbf{x} \in \mathcal{X}} |R(\mathbf{x}) - R_n(\mathbf{x})|$$

○ Main tool: concentration inequalities!

▶ Measure of how far is an empirical average from the true mean

## Theorem (Hoeffding's Inequality [9])

*Let $Y_1, \ldots, Y_n$ be i.i.d. random variables with $Y_i$ taking values in the interval $[a_i, b_i] \subseteq \mathbb{R}$ for all $i = 1, \ldots, n$. Let $S_n := \frac{1}{n} \sum_{i=1}^{n} Y_i$. It holds that*

$$\mathbb{P}\left(|S_n - \mathbb{E}[S_n]| > t\right) \leq 2 \exp\left(-\frac{2n^2 t^2}{\sum_{i=1}^{n}(b_i - a_i)^2}\right)$$

## Warmup: Generalization bound for a singleton

### Lemma

*For $i = 1, \ldots, n$, let $(\mathbf{a}_i, b_i) \in \mathbb{R}^p \times \{-1, 1\}$ be independent random variables and $h_{\mathbf{x}} : \mathbb{R}^p \to \mathbb{R}$ be a function parametrized by $\mathbf{x} \in \mathcal{X}$. Let $\mathcal{X} = \{\mathbf{x}_0\}$ and $L(h_{\mathbf{x}}(\mathbf{a}), b) = \{sign(h_{\mathbf{x}}(\mathbf{a})) \neq b\}$ be the 0-1 loss. With probability at least $1 - \delta$, we have that*

$$\sup_{\mathbf{x} \in \mathcal{X}} |R(\mathbf{x}) - R_n(\mathbf{x})| = |R(\mathbf{x}_0) - R_n(\mathbf{x}_0)| \leq \sqrt{\frac{\ln(2/\delta)}{2n}}.$$

### Proof.

Note that $\mathbb{E}[\frac{1}{n} \sum_{i=1}^{n} L(h_{\mathbf{x}_0}(\mathbf{a}_i), b_i)] = R(\mathbf{x}_0)$, the expected risk of the parameter $\mathbf{x}_0$. Moreover $L(h_{\mathbf{x}_0}(\mathbf{a}_i), b_i) \in [0, 1]$. We can use Hoeffding's inequality and obtain

$$\mathbb{P}(|R_n(\mathbf{x}_0) - R(\mathbf{x}_0)| > t) = \mathbb{P}\left( \left| \frac{1}{n} \sum_{i=1}^{n} L_i(h_{\mathbf{x}_0}(\mathbf{a}_i), b_i) - R(\mathbf{x}_0) \right| > t \right) \leq 2 \exp\left( -2nt^2 \right)$$

Setting $\delta := 2 \exp\left( -2nt^2 \right)$ we have that $t = \sqrt{\frac{\ln 2/\delta}{2n}}$, thus obtaining the result.

$\square$

## Generalization bound for finite sets

### Lemma

*For $i = 1, \ldots, n$, let $(\mathbf{a}_i, b_i) \in \mathbb{R}^p \times \{-1, 1\}$ be independent random variables and $h_{\mathbf{x}} : \mathbb{R}^p \to \mathbb{R}$ be a function parametrized by $\mathbf{x} \in \mathcal{X}$. Let $\mathcal{X}$ be a finite set and $L(h_{\mathbf{x}}(\mathbf{a}), b) = \{sign(h_{\mathbf{x}}(\mathbf{a})) \neq b\}$ be the 0-1 loss. With probability at least $1 - \delta$, we have that*

$$\sup_{\mathbf{x} \in \mathcal{X}} |R(\mathbf{x}) - R_n(\mathbf{x})| \leq \sqrt{\frac{\ln |\mathcal{X}| + \ln(2/\delta)}{2n}}.$$

### Proof.

Let $\mathcal{X} = \{\mathbf{x}_1, \ldots, \mathbf{x}_{|\mathcal{X}|}\}$. We can use a union bound and the analysis of the singleton case to obtain:

$$\mathbb{P}(\exists j : |R_n(\mathbf{x}_j) - R(\mathbf{x}_j)| > t) \leq \sum_{j=1}^{|\mathcal{X}|} \mathbb{P}(|R_n(\mathbf{x}_j) - R(\mathbf{x}_j)| > t) = 2|\mathcal{X}| \exp\left(-2nt^2\right)$$

Setting $\delta := 2|\mathcal{X}| \exp\left(-2nt^2\right)$, we have that $t = \sqrt{\frac{\ln |\mathcal{X}| + \ln \frac{2}{\delta}}{2n}}$, thus obtaining the result.

□

# Generalization bounds for infinite classes - The Rademacher complexity

<span style="color:red">However, in most applications in ML/DL we optimize over an infinite parameter space $\mathcal{X}$!</span>

○ Need a notion of *complexity* to derive generalization bounds for infinite classes of functions

---

### Definition (Rademacher Complexity [3])

Let $S = \{\mathbf{a}_1, \ldots, \mathbf{a}_n\} \subseteq \mathbb{R}^p$ and let $\{\sigma_i : i = 1, \ldots, n\}$ be independent Rademacher random variables i.e., taking values uniformly in $\{-1, +1\}$ (coin flip). Let $\mathcal{H}$ be a class of functions of the form $h : \mathbb{R}^p \to \mathbb{R}$. The Rademacher complexity of $\mathcal{H}$ <span style="color:red">with respect to $A$</span> is defined as:

$$\mathcal{R}_A(\mathcal{H}) := \mathbb{E} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_i h(\mathbf{a}_i).$$

---

$\mathcal{R}_A(\mathcal{H})$ measures how well can we fit random signs ($\pm 1$) with the output of an element of $\mathcal{H}$ on the set $A$.

# Visualizing Rademacher complexity
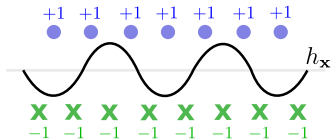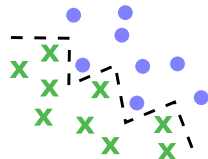


Figure: Rademacher complexity measures correlation with random signs

# Visualizing Rademacher complexity



(a) High Rademacher Complexity

(b) Large Generalization error (memorization)

(c) Low Rademacher Complexity

(d) Low Generalization error

Figure: Rademacher complexity and Generalization error

# Fundamental theorem about the Rademacher Complexity

## Theorem (See Theorem 3.3 and 5.8 in [9])

*Suppose that the loss function has the form $L(h_\mathbf{x}(\mathbf{a}), b) = \phi(b \cdot h_\mathbf{x}(\mathbf{a}))$ for a 1-Lipschitz function $\phi : \mathbb{R} \to \mathbb{R}$.*

*Let $\mathcal{H}_\mathcal{X} := \{h_\mathbf{x} : \mathbf{x} \in \mathcal{X}\}$ be a class of parametric functions $h_\mathbf{x} : \mathbb{R}^p \to \mathbb{R}$. For any $\delta > 0$, with probability at least $1 - \delta$ over the draw of an i.i.d. sample $\{(\mathbf{a}_i, b_i)\}_{i=1}^n$, letting $A = (\mathbf{a}_1, \ldots, \mathbf{a}_n)$, the following holds:*

$$\sup_{\mathbf{x} \in \mathcal{X}} |R_n(\mathbf{x}) - R(\mathbf{x})| \leq 2\mathbb{E}_A \mathcal{R}_A(\mathcal{H}_\mathcal{X}) + \sqrt{\frac{\ln(2/\delta)}{2n}}$$

$$\sup_{\mathbf{x} \in \mathcal{X}} |R_n(\mathbf{x}) - R(\mathbf{x})| \leq 2\mathcal{R}_A(\mathcal{H}_\mathcal{X}) + 3\sqrt{\frac{\ln(4/\delta)}{2n}}$$

## Assumption is true for common losses

- $L(h_\mathbf{x}(\mathbf{a}), b) = \log(1 + \exp(-b \cdot h_\mathbf{x}(\mathbf{a}))) \Rightarrow \phi(z) := \log(1 + \exp(z))$ (logistic loss)
- $L(h_\mathbf{x}(\mathbf{a}), b) = \max(0, 1 - b \cdot h_\mathbf{x}(\mathbf{a})) \Rightarrow \phi(z) := \max(0, 1 - z)$ (hinge loss)

# Computing the Rademacher complexity of linear functions

## Theorem

Let $\mathcal{X} := \{\mathbf{x} \in \mathbb{R}^p : \|\mathbf{x}\|_2 \leq \lambda\}$ and let $\mathcal{H}_\mathcal{X}$ be the class of functions of the form $h_\mathbf{x} : \mathbb{R}^p \to \mathbb{R}, h_\mathbf{x}(\mathbf{a}) = \langle \mathbf{x}, \mathbf{a} \rangle$, for some $\mathbf{x} \in \mathcal{X}\}$. Let $A = \{\mathbf{a}_1, \ldots, \mathbf{a}_n\} \subseteq \mathbb{R}^p$ such that $\max_{i=1,\ldots,n} \|\mathbf{a}_i\| \leq M$. It holds that $\mathcal{R}_A(\mathcal{H}_\mathcal{X}) \leq \lambda M / \sqrt{n}$.

## Proof.

$$\mathcal{R}_A(\mathcal{H}_\mathcal{X}) = \mathbb{E} \sup_{\|\mathbf{x}\|_2 \leq \lambda} \frac{1}{n} \sum_{i=1}^{n} \sigma_i \langle \mathbf{x}, \mathbf{a} \rangle$$

$$= \mathbb{E} \sup_{\|\mathbf{x}\|_2 \leq \lambda} \frac{1}{n} \left\langle \mathbf{x}, \sum_{i=1}^{n} \sigma_i \mathbf{a} \right\rangle$$

$$\leq \frac{1}{n} \lambda \mathbb{E} \left\| \sum_{i=1}^{n} \sigma_i \mathbf{a}_i \right\|_2 \quad \text{(C-S)}$$

$$\Rightarrow \mathcal{R}_A(\mathcal{H}_\mathcal{X}) \leq \frac{1}{n} \lambda \left( \mathbb{E} \sum_{i=1}^{n} \|\sigma_i \mathbf{a}_i\|_2^2 \right)^{1/2} \quad \text{(Jensen)}$$

$$\leq \frac{1}{n} \lambda \left( \sum_{i=1}^{n} \|\mathbf{a}_i\|_2^2 \right)^{1/2}$$

$$\leq \lambda M / \sqrt{n}$$

$\square$

## Rademacher complexity estimates of fully connected Neural Networks

### Notation

For a matrix $\mathbf{X} \in \mathbb{R}^{n,m}$, $\|\mathbf{X}\|$ denotes its spectral norm. Let $\mathbf{X}_{:,k}$ be the $k$-th column of $\mathbf{X}$. We define

$$\|\mathbf{X}\|_{2,1} = \|(\|\mathbf{X}_{:,1}\|_2, \ldots, \|\mathbf{X}_{:,m}\|_2)\|_1. \tag{2}$$

### Theorem (Spectral bound [2])

*For positive integers $p_0, p_1, \ldots, p_d = 1$, and positive reals $\lambda_1, \ldots, \lambda_d$ and $\nu_1, \ldots, \nu_d$, define the set*

$$\mathcal{X} := \{(\mathbf{X}_1, \ldots, \mathbf{X}_d) : \mathbf{X}_i \in \mathbb{R}^{p_i \times p_{i-1}}, \|\mathbf{X}_i\| \leq \lambda_i, \|\mathbf{X}_i^T\|_{2,1} \leq \nu_i\}$$

*Let $H_{\mathcal{X}}$ be the class of neural networks $h_{\mathbf{x}} : \mathbb{R}^p \to \mathbb{R}$, $h_{\mathbf{x}} = \mathbf{X}_d \circ \sigma \circ \ldots \circ \sigma \circ \mathbf{X}_1$ where $\mathbf{x} = (\mathbf{X}_1, \ldots, \mathbf{X}_d) \in \mathcal{X}$. Suppose that $\sigma$ is 1-Lipschitz. Let $A = \{\mathbf{a}_1, \ldots, \mathbf{a}_n\} \subseteq \mathbb{R}^p$, $M := \max_{i=1,\ldots,n} \|\mathbf{a}_i\|$ and $W := \max\{p_i : i = 0, \ldots, d\}$.*

*The Rademacher complexity of $\mathcal{H}_{\mathcal{X}}$ with respect to $A$ is bounded as*

$$\mathcal{R}_A(\mathcal{H}_{\mathcal{X}}) = \mathcal{O}\left(\frac{\log(W)M}{\sqrt{n}} \prod_{i=1}^{d} \lambda_i \left(\sum_{j=1}^{d} \frac{\nu_j^{2/3}}{\lambda_j^{2/3}}\right)^{3/2}\right). \tag{3}$$

# How well do complexity measures correlate with generalization?

| name | definition | correlation[1] |
|---|---|---|
| Frobenius distance to initialization [10] | $\sum_{i=1}^{d} \|\mathbf{X}_i - \mathbf{X}_i^0\|_F^2$ | $-0.263$ |
| Spectral complexity[2] [2] | $\prod_{i=1}^{d} \|\mathbf{X}_i\| \left( \sum_{i=1}^{d} \frac{\|\mathbf{X}_i\|_{2,1}^{3/2}}{\|\mathbf{X}_i\|^{3/2}} \right)^{2/3}$ | $\mathbf{-0.537}$ |
| Parameter Frobenius norm | $\sum_{i=1}^{d} \|\mathbf{X}_i\|_F^2$ | $0.073$ |
| Fisher-Rao [8] | $\frac{(d+1)^2}{n} \sum_{i=1}^{n} \langle \mathbf{x}, \nabla_{\mathbf{x}} \ell(h_{\mathbf{x}}(\mathbf{a}_i), b_i) \rangle$ | $0.078$ |
| Path-norm [11] | $\sum_{(i_0,\ldots,i_d)} \prod_{j=1}^{d} \left( \mathbf{X}_{i_j, i_{j-1}} \right)^2$ | $\mathbf{0.373}$ |

Table: Complexity measures compared in the empirical study [7], and their correlation with generalization

Complexity measures are still far from explaining generalization in Deep Learning!

A more recent evaluation of many complexity measures is available [5].

---

[1] Kendall's rank correlation coefficient
[2] The definition in [7] differs slightly

# Wrap up!

○ Deep learning tricks-of-the-trade recitation on Friday!

$$R(\mathbf{x}^t) - R(\mathbf{x}^\natural) = R(\mathbf{x}^t) - R_n(\mathbf{x}^t) + R_n(\mathbf{x}^t) - R_n(\mathbf{x}^\star) + \underbrace{R_n(\mathbf{x}^\star) - R_n(\mathbf{x}^\natural)}_{\leq 0} + R_n(\mathbf{x}^\natural) - R(\mathbf{x}^\natural)$$

$$\leq R_n(\mathbf{x}^t) - R_n(\mathbf{x}^\star) + \underbrace{R(\mathbf{x}^t) - R_n(\mathbf{x}^t) + R_n(\mathbf{x}^\natural) - R(\mathbf{x}^\natural)}_{2\sup_{\mathbf{x}\in\mathcal{X}}|R_n(\mathbf{x})-R(\mathbf{x})|}$$

$$R(\mathbf{x}^t) - R(\mathbf{x}^\circ) = R(\mathbf{x}^t) - R(\mathbf{x}^\natural) + R(\mathbf{x}^\natural) - R(\mathbf{x}^\circ)$$

$$\leq R_n(\mathbf{x}^t) - R_n(\mathbf{x}^\star) + 2\sup_{\mathbf{x}\in\mathcal{X}}|R_n(\mathbf{x}) - R(\mathbf{x})| + R(\mathbf{x}^\natural) - R(\mathbf{x}^\circ)$$

# References I

[1] Anish Athalye, Nicholas Carlini, and David Wagner.
Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples.
In *International conference on machine learning*, pages 274–283. PMLR, 2018.

[2] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky.
Spectrally-normalized margin bounds for neural networks.
In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 6240–6249. Curran Associates, Inc., 2017.

[3] Peter L Bartlett and Shahar Mendelson.
Rademacher and gaussian complexities: Risk bounds and structural results.
*Journal of Machine Learning Research*, 3(Nov):463–482, 2002.

[4] George Cybenko.
Approximation by superpositions of a sigmoidal function.
*Mathematics of control, signals and systems*, 2(4):303–314, 1989.

[5] Gintare Karolina Dziugaite, Alexandre Drouin, Brady Neal, Nitarshan Rajkumar, Ethan Caballero, Linbo Wang, Ioannis Mitliagkas, and Daniel M Roy.
In search of robust measures of generalization.
*arXiv preprint arXiv:2010.11924*, 2020.

# References II

[6] Kevin Eykholt, Ivan Evtimov, Earlence Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Song.
Robust physical-world attacks on deep learning visual classification.
In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1625–1634, 2018.

[7] Yiding Jiang*, Behnam Neyshabur*, Hossein Mobahi, Dilip Krishnan, and Samy Bengio.
Fantastic generalization measures and where to find them.
In *International Conference on Learning Representations*, 2020.

[8] Tengyuan Liang, Tomaso Poggio, Alexander Rakhlin, and James Stokes.
Fisher-rao metric, geometry, and complexity of neural networks.
volume 89 of *Proceedings of Machine Learning Research*, pages 888–896. PMLR, 16–18 Apr 2019.

[9] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar.
*Foundations of Machine Learning*.
The MIT Press, 2nd edition, 2018.

[10] Vaishnavh Nagarajan and J. Zico Kolter.
Generalization in Deep Networks: The Role of Distance from Initialization.
*arXiv e-prints*, page arXiv:1901.01672, January 2019.

# References III

[11] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro.
Norm-based capacity control in neural networks.
In *Conference on Learning Theory*, pages 1376–1401, 2015.

[12] Ali Shafahi, Mahyar Najibi, Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis,
Gavin Taylor, and Tom Goldstein.
Adversarial training for free!
*arXiv preprint arXiv:1904.12843*, 2019.