

The Mathematics of Large Machine Learning Models

Andrea Montanari

Stanford University

August 11, 2025

Prelude: The role of theory

Current state of Machine Learning

I started working in high-dim statistics and ML and around 2008

- ▶ **2008:** $\approx 1.5 \cdot 10^3$ papers submitted each year in top 3 venues
- ▶ **2024:** $\approx 3 \cdot 10^4$ papers submitted each year in top 3 venues
- ▶ Major discoveries
 - ▶ Machines can learn language
 - ▶ Learn from examples during conversations
 - ▶ Exhibit step-by-step problem-solving

Current state of Machine Learning

I started working in high-dim statistics and ML and around 2008

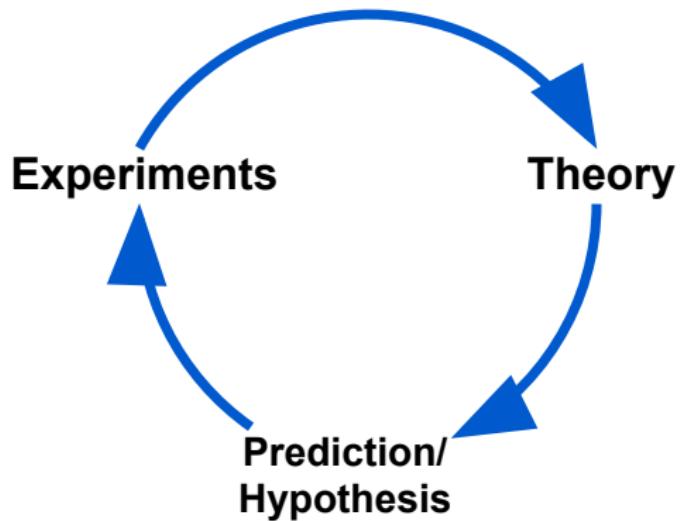
- ▶ **2008:** $\approx 1.5 \cdot 10^3$ papers submitted each year in top 3 venues
- ▶ **2024:** $\approx 3 \cdot 10^4$ papers submitted each year in top 3 venues
- ▶ Major discoveries
 - ▶ Machines can learn language
 - ▶ Learn from examples during conversations
 - ▶ Exhibit step-by-step problem-solving

Current state of Machine Learning

I started working in high-dim statistics and ML and around 2008

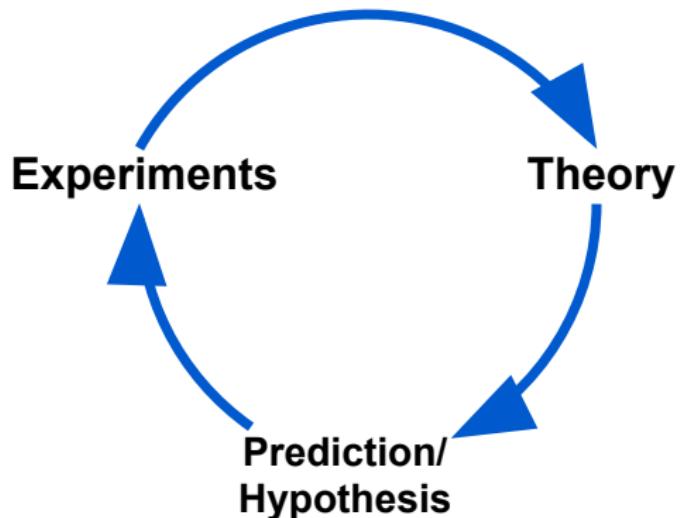
- ▶ **2008:** $\approx 1.5 \cdot 10^3$ papers submitted each year in top 3 venues
- ▶ **2024:** $\approx 3 \cdot 10^4$ papers submitted each year in top 3 venues
- ▶ Major discoveries
 - ▶ Machines can learn language
 - ▶ Learn from examples during conversations
 - ▶ Exhibit step-by-step problem-solving

A success of the scientific method?



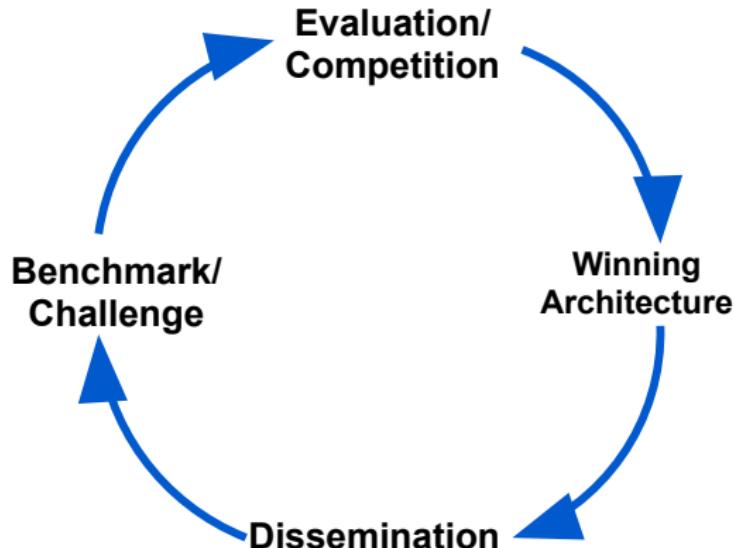
Not really!

A success of the scientific method?

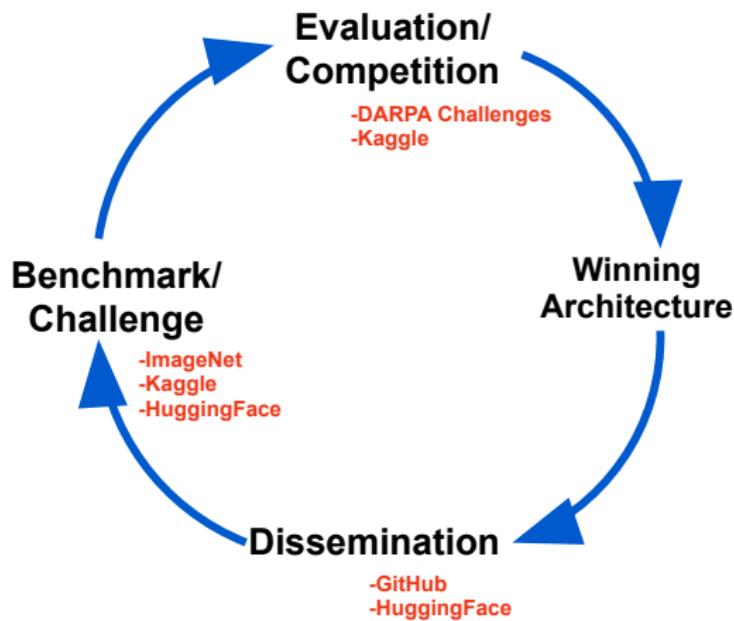


Not really!

The 'AI research method'



The AI research method and its institutions

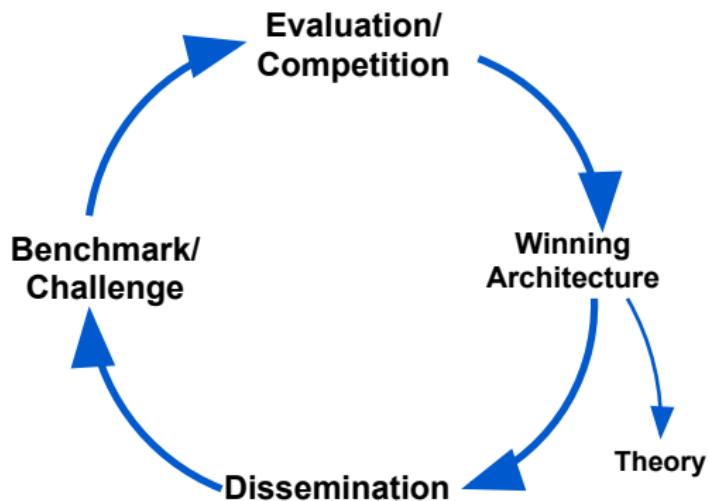


The AI research method and its institutions



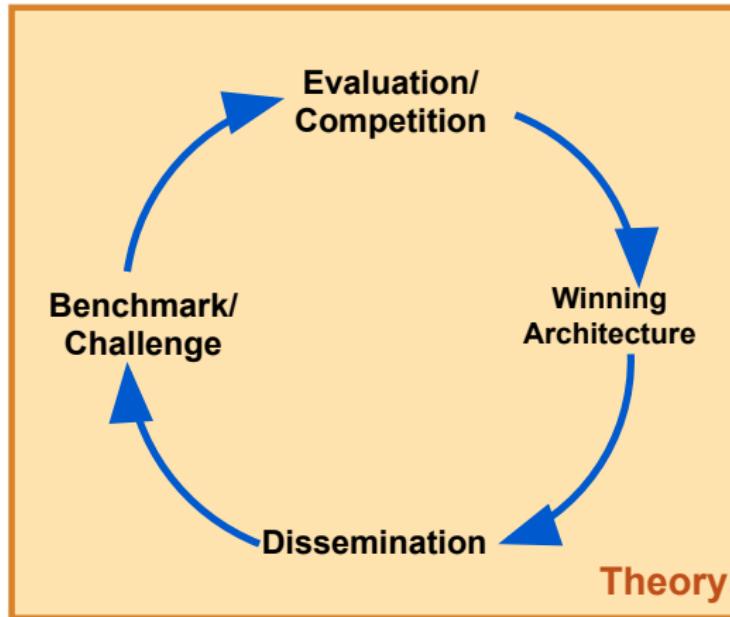
A screenshot of a GitHub repository page for 'karpathy / nanoGPT'. The page includes a header with navigation links for Code, Issues, Pull requests, Actions, Projects, and Security. Below the header, it shows the repository name 'nanoGPT' and its status as 'Public'. It displays metrics like 222 issues, 67 pull requests, 40.6k stars, and 6.7k forks. The main content area shows the repository's structure with files like assets, config, data, .gitattributes, .gitignore, and LICENSE. To the right, there is an 'About' section with a brief description: 'The simplest, fastest repository for training/finetuning medium-sized GPTs.' and links to Readme, MIT license, Activity, 40.6k stars, 401 watching, 6.7k forks, and Report repository.

And theory?



- ▶ Explain why architecture X is currently winning the race

And theory?



- ▶ Ask very fundamental questions about very general phenomena
 - ▶ Does more data always help? How much?
 - ▶ Do more complex models always help? How much?

A few theoretical successes

- ▶ Overfitting and generalization
- ▶ Scaling
- ▶ Generic phenomena in (stochastic) gradient descent
- ▶ Benchmarking/uncertainty quantification
- ▶ ...

A few theoretical successes

- ▶ **Overfitting and generalization**
- ▶ Scaling
- ▶ Generic phenomena in (stochastic) gradient descent
- ▶ Benchmarking/uncertainty quantification
- ▶ ...

Outline

- ➊ The generalization problem
- ➋ Enters modern machine learning
- ➌ Two previews
- ➍ Conclusion

The generalization problem

History

Before 1960:

- ▶ A program defines the input-output relation

Rosenblatt, 1960:

- ▶ Input-output pairs (x_i, y_i) , $i \leq n$
- ▶ Let the machine **learn** the relation
- ▶ Proof of concept: ‘Perceptron’

History

Before 1960:

- ▶ A program defines the input-output relation

Rosenblatt, 1960:

- ▶ Input-output pairs (x_i, y_i) , $i \leq n$
- ▶ Let the machine **learn** the relation
- ▶ Proof of concept: ‘Perceptron’

History

Before 1960:

- ▶ A program defines the input-output relation

Rosenblatt, 1960:

- ▶ Input-output pairs (x_i, y_i) , $i \leq n$
- ▶ Let the machine **learn** the relation
- ▶ Proof of concept: ‘Perceptron’



Frank Rosenblatt

FIG. 1 — Organization of a biological brain. (Red areas indicate active cells, responding to the letter X.)

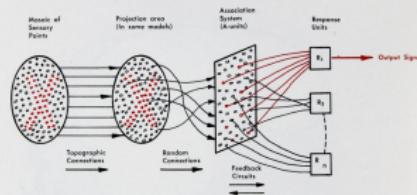


FIG. 2 — Organization of a perceptron.

Perceptron

History

Before 1960:

- ▶ A program defines the input-output relation

Rosenblatt, 1960:

- ▶ Input-output pairs (x_i, y_i) , $i \leq n$
- ▶ Let the machine **learn** the relation
- ▶ Proof of concept: ‘Perceptron’

Why does it work? What does it mean to learn?

History

Before 1960:

- ▶ A computer is programmed
- ▶ A program needs to determine the input-output relation

Rosenblatt 1960:

- ▶ Show the machine input-output pairs (x_i, y_i) , $i \leq n$
- ▶ Let the machine **learn** the relation
- ▶ Proof of concept: Perceptron

Vapnik–Chervonenkis, 1969-1974:

- ▶ Learning is statistical learning

History

Before 1960:

- ▶ A computer is programmed
- ▶ A program needs to determine the input-output relation

Rosenblatt 1960:

- ▶ Show the machine input-output pairs (x_i, y_i) , $i \leq n$
- ▶ Let the machine **learn** the relation
- ▶ Proof of concept: Perceptron

Vapnik–Chervonenkis, 1969-1974:

- ▶ Learning is statistical learning

Vapnik studied in Samarkand: closer to Bangalore than Stanford is to MIT.

Statistical learning

- ▶ **Data:** $(x_1, y_1), \dots, (x_n, y_n) \sim_{\text{iid}} \mathbb{P}$
- ▶ **Input-output relations:** $f(\cdot; \theta_1), f(\cdot; \theta_2), \dots, f(\cdot; \theta_M)$.

Statistical learning

- ▶ **Data:** $(x_1, y_1), \dots, (x_n, y_n) \sim_{\text{iid}} \mathbb{P}$
- ▶ **Models:** $f(\cdot; \theta_1), f(\cdot; \theta_2), \dots, f(\cdot; \theta_M)$.
- ▶ **Learning:** $\mathcal{R}(\theta) := \mathbb{E}\{\text{dist}(y_{n+1}, f(x_{n+1}; \theta))\}$:
 - ▶ $\hat{\theta}$ selected from data $(x_1, y_1), \dots, (x_n, y_n)$
 - ▶ Minimize test error $\mathcal{R}(\theta)$!

Statistical learning

- ▶ **Data:** $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \sim_{\text{iid}} \mathbb{P}$
- ▶ **Models:** $f(\cdot; \boldsymbol{\theta}_1), f(\cdot; \boldsymbol{\theta}_2), \dots, f(\cdot; \boldsymbol{\theta}_M)$.
- ▶ **Learning:** $\mathcal{R}(\boldsymbol{\theta}) := \mathbb{E}\{\text{dist}(y_{n+1}, f(\mathbf{x}_{n+1}; \boldsymbol{\theta}))\}$:
 - ▶ $\hat{\boldsymbol{\theta}}$ selected from data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
 - ▶ Minimize test error $\mathcal{R}(\boldsymbol{\theta})$!

Statistical learning

- ▶ **Data:** $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \sim_{\text{iid}} \mathbb{P}$
- ▶ **Models:** $f(\cdot; \boldsymbol{\theta}_1), f(\cdot; \boldsymbol{\theta}_2), \dots, f(\cdot; \boldsymbol{\theta}_M)$.
- ▶ **Learning:** $\mathcal{R}(\boldsymbol{\theta}) := \mathbb{E}\{\text{dist}(y_{n+1}, f(\mathbf{x}_{n+1}; \boldsymbol{\theta}))\}$:
 - ▶ $\hat{\boldsymbol{\theta}}$ selected from data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
 - ▶ Minimize test error $\mathcal{R}(\boldsymbol{\theta})$!

Statisticians: ‘This is *just* statistics’

Statistical learning

$\text{dist}(y_i, f(x_i; \theta))$	$f(\cdot; \theta_1)$	$f(\cdot; \theta_2)$	\dots	$f(\cdot; \theta_M)$
(y_1, x_1)	0	1	\dots	0
(y_2, x_2)	1	0	\dots	1
\vdots	\vdots	\vdots	\ddots	\vdots
(y_n, x_n)	0	1	\dots	1
(y_{n+1}, x_{n+1})	1	0	\dots	1

	$f(\cdot; \theta_1)$	$f(\cdot; \theta_2)$	\dots	$f(\cdot; \theta_M)$
empirical risk	$\hat{\mathcal{R}}_n(\theta_1)$	$\hat{\mathcal{R}}_n(\theta_2)$	\dots	$\hat{\mathcal{R}}_n(\theta_M)$
population risk	$\mathcal{R}(\theta_1)$	$\mathcal{R}(\theta_2)$	\dots	$\mathcal{R}(\theta_M)$

$$\hat{\mathcal{R}}_n(\theta) := \frac{1}{n} \sum_{i=1}^n \text{dist}(y_i, f(x_i; \theta)), \quad \mathcal{R}(\theta) = \mathbb{E} \text{dist}(y_{n+1}, f(x_{n+1}; \theta)).$$

Statistical learning

$\text{dist}(y_i, f(x_i; \theta))$	$f(\cdot; \theta_1)$	$f(\cdot; \theta_2)$	\dots	$f(\cdot; \theta_M)$
(y_1, x_1)	0	1	\dots	0
(y_2, x_2)	1	0	\dots	1
\vdots	\vdots	\vdots	\ddots	\vdots
(y_n, x_n)	0	1	\dots	1
(y_{n+1}, x_{n+1})	1	0	\dots	1

	$f(\cdot; \theta_1)$	$f(\cdot; \theta_2)$	\dots	$f(\cdot; \theta_M)$
empirical risk	$\hat{\mathcal{R}}_n(\theta_1)$	$\hat{\mathcal{R}}_n(\theta_2)$	\dots	$\hat{\mathcal{R}}_n(\theta_M)$
population risk	$\mathcal{R}(\theta_1)$	$\mathcal{R}(\theta_2)$	\dots	$\mathcal{R}(\theta_M)$

$$\hat{\mathcal{R}}_n(\theta) := \frac{1}{n} \sum_{i=1}^n \text{dist}(y_i, f(x_i; \theta)), \quad \mathcal{R}(\theta) = \mathbb{E} \text{dist}(y_{n+1}, f(x_{n+1}; \theta)).$$

Statistical learning

$\text{dist}(y_i, f(x_i; \theta))$	$f(\cdot; \theta_1)$	$f(\cdot; \theta_2)$	\dots	$f(\cdot; \theta_M)$
(y_1, x_1)	0	1	\dots	0
(y_2, x_2)	1	0	\dots	1
\vdots	\vdots	\vdots	\ddots	\vdots
(y_n, x_n)	0	1	\dots	1
(y_{n+1}, x_{n+1})	1	0	\dots	1

	$f(\cdot; \theta_1)$	$f(\cdot; \theta_2)$	\dots	$f(\cdot; \theta_M)$
empirical risk	$\hat{\mathcal{R}}_n(\theta_1)$	$\hat{\mathcal{R}}_n(\theta_2)$	\dots	$\hat{\mathcal{R}}_n(\theta_M)$
population risk	$\mathcal{R}(\theta_1)$	$\mathcal{R}(\theta_2)$	\dots	$\mathcal{R}(\theta_M)$

$$\hat{\mathcal{R}}_n(\theta) := \frac{1}{n} \sum_{i=1}^n \text{dist}(y_i, f(x_i; \theta)), \quad \mathcal{R}(\theta) = \mathbb{E} \text{dist}(y_{n+1}, f(x_{n+1}; \theta)).$$

Statistical learning

$$\widehat{\mathcal{R}}_n(\theta) := \frac{1}{n} \sum_{i=1}^n \text{dist}(y_i, f(x_i)), \quad \mathcal{R}(\theta) = \mathbb{E} \text{dist}(y_i, f(x_i; \theta)).$$

- Central limit theorem: $|\mathcal{R}(\theta_i) - \widehat{\mathcal{R}}_n(\theta_i)| = O(1/\sqrt{n})$
- Vapnik–Chervonenkis, 1969-1974: $\mathcal{F} := \{f(\cdot; \theta) : \theta \in \mathcal{S}\}$

$$\max_{\theta \in \mathcal{S}} |\widehat{\mathcal{R}}_n(\theta) - \mathcal{R}(\theta)| \leq \text{const.} \sqrt{\frac{\text{dvc}(\mathcal{F})}{n}}$$

Statistical learning

$$\widehat{\mathcal{R}}_n(\theta) := \frac{1}{n} \sum_{i=1}^n \text{dist}(y_i, f(x_i)), \quad \mathcal{R}(\theta) = \mathbb{E} \text{dist}(y_i, f(x_i; \theta)).$$

- Central limit theorem: $|\mathcal{R}(\theta_i) - \widehat{\mathcal{R}}_n(\theta_i)| = O(1/\sqrt{n})$
- Vapnik–Chervonenkis, 1969-1974: $\mathcal{F} := \{f(\cdot; \theta) : \theta \in \mathcal{S}\}$

$$\max_{\theta \in \mathcal{S}} |\widehat{\mathcal{R}}_n(\theta) - \mathcal{R}(\theta)| \leq \text{const.} \sqrt{\frac{\text{dvc}(\mathcal{F})}{n}}$$

Statistical learning

$$\widehat{\mathcal{R}}_n(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \text{dist}(y_i, f(\mathbf{x}_i)), \quad \mathcal{R}(\boldsymbol{\theta}) = \mathbb{E} \text{dist}(y_i, f(\mathbf{x}_i; \boldsymbol{\theta})).$$

- Central limit theorem: $|\mathcal{R}(\boldsymbol{\theta}_i) - \widehat{\mathcal{R}}_n(\boldsymbol{\theta}_i)| = O(1/\sqrt{n})$
- Vapnik–Chervonenkis, 1969-1974: $\mathcal{F} := \{f(\cdot; \boldsymbol{\theta}) : \boldsymbol{\theta} \in \mathcal{S}\}$

$$\max_{\boldsymbol{\theta} \in \mathcal{S}} |\widehat{\mathcal{R}}_n(\boldsymbol{\theta}) - \mathcal{R}(\boldsymbol{\theta})| \leq \text{const.} \sqrt{\frac{\text{dvc}(\mathcal{F})}{n}}$$

Statistical learning

$$\widehat{\mathcal{R}}_n(\theta) := \frac{1}{n} \sum_{i=1}^n \text{dist}(y_i, f(x_i)), \quad \mathcal{R}(\theta) = \mathbb{E} \text{dist}(y_i, f(x_i; \theta)).$$

- Central limit theorem: $|\mathcal{R}(\theta_i) - \widehat{\mathcal{R}}_n(\theta_i)| = O(1/\sqrt{n})$
- Vapnik–Chervonenkis, 1969-1974: $\mathcal{F} := \{f(\cdot; \theta) : \theta \in \mathcal{S}\}$

$$\max_{\theta \in \mathcal{S}} |\widehat{\mathcal{R}}_n(\theta) - \mathcal{R}(\theta)| \lesssim \sqrt{\frac{\text{dvc}(\mathcal{F})}{n}}$$

Key insight: Uniform CLT-style bounds

$$\max_{\theta \in \mathcal{S}} |\hat{\mathcal{R}}_n(\theta) - \mathcal{R}(\theta)| \lesssim \sqrt{\frac{dvc(\mathcal{F})}{n}}$$

- ▶ Simultaneously for infinitely many θ 's!
- ▶ Bound does not depend on the number of parameters p ($\theta \in \mathbb{R}^p$)
- ▶ Only depends on the function class \mathcal{F} via $dvc(\mathcal{F})$.

Key insight: Uniform CLT-style bounds

$$\max_{f \in \mathcal{F}} |\hat{\mathcal{R}}_n(f) - \mathcal{R}(f)| \lesssim \sqrt{\frac{d_{VC}(\mathcal{F})}{n}}$$

- ▶ Simultaneously for infinitely many θ 's!
- ▶ Bound does not depend on the number of parameters p ($\theta \in \mathbb{R}^p$)
- ▶ Only depends on the function class \mathcal{F} via $d_{VC}(\mathcal{F})$.

A couple of examples

Example 1: Linear model

$$\mathcal{F}_{\text{lin}}(\rho) := \left\{ f(x; \theta) = \langle \theta, x \rangle : \|\theta\| \leq \rho \right\},$$

$$\max_{f \in \mathcal{F}_{\text{lin}}(\rho)} |\widehat{\mathcal{R}}_n(f) - \mathcal{R}(f)| \lesssim \rho \sqrt{\frac{d}{n}}.$$

Example 2: 2-layer networks

[Bartlett 1996]

$$\mathcal{F}_{\text{2lnn}}(\gamma) := \left\{ f(x; \theta) = \sum_{i=1}^N a_i \sigma(\langle w_i, x \rangle) : \|w_i\| = 1 \forall i \in [N], \|a\|_1 \leq \gamma \right\},$$

$$\max_{f \in \mathcal{F}_{\text{2lnn}}(\gamma)} |\widehat{\mathcal{R}}_n(f) - \mathcal{R}(f)| \lesssim \gamma \sqrt{\frac{d}{n}}. \quad \text{Independent of } N!$$

A couple of examples

Example 1: Linear model

$$\mathcal{F}_{\text{lin}}(\rho) := \left\{ f(x; \theta) = \langle \theta, x \rangle : \|\theta\| \leq \rho \right\},$$

$$\max_{f \in \mathcal{F}_{\text{lin}}(\rho)} |\widehat{\mathcal{R}}_n(f) - \mathcal{R}(f)| \lesssim \rho \sqrt{\frac{d}{n}}.$$

Example 2: 2-layer networks

[Bartlett 1996]

$$\mathcal{F}_{\text{2lnn}}(\gamma) := \left\{ f(x; \theta) = \sum_{i=1}^N a_i \sigma(\langle w_i, x \rangle) : \|w_i\| = 1 \forall i \in [N], \|a\|_1 \leq \gamma \right\},$$

$$\max_{f \in \mathcal{F}_{\text{2lnn}}(\gamma)} |\widehat{\mathcal{R}}_n(f) - \mathcal{R}(f)| \lesssim \gamma \sqrt{\frac{d}{n}}. \quad \text{Independent of } N!$$

A couple of examples

Example 1: Linear model

$$\mathcal{F}_{\text{lin}}(\rho) := \left\{ f(x; \theta) = \langle \theta, x \rangle : \|\theta\| \leq \rho \right\},$$

$$\max_{f \in \mathcal{F}_{\text{lin}}(\rho)} |\widehat{\mathcal{R}}_n(f) - \mathcal{R}(f)| \lesssim \rho \sqrt{\frac{d}{n}}.$$

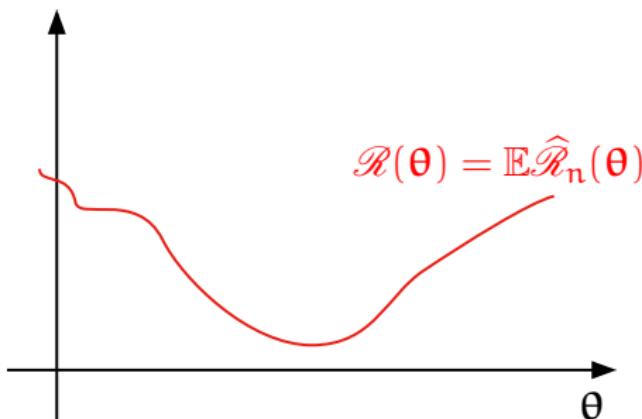
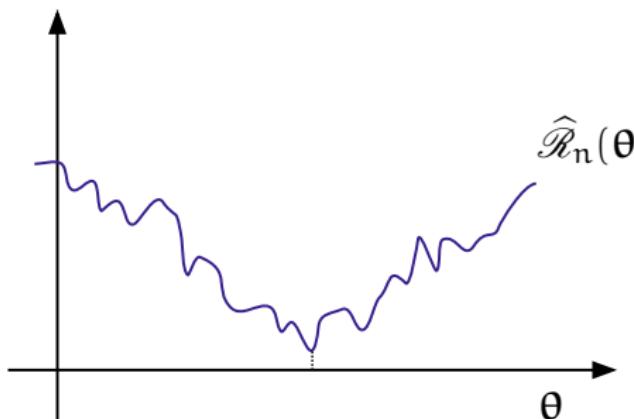
Example 2: 2-layer networks

[Bartlett 1996]

$$\mathcal{F}_{\text{2lnn}}(\gamma) := \left\{ f(x; \theta) = \sum_{i=1}^N a_i \sigma(\langle w_i, x \rangle) : \|w_i\| = 1 \forall i \in [N], \|a\|_1 \leq \gamma \right\},$$

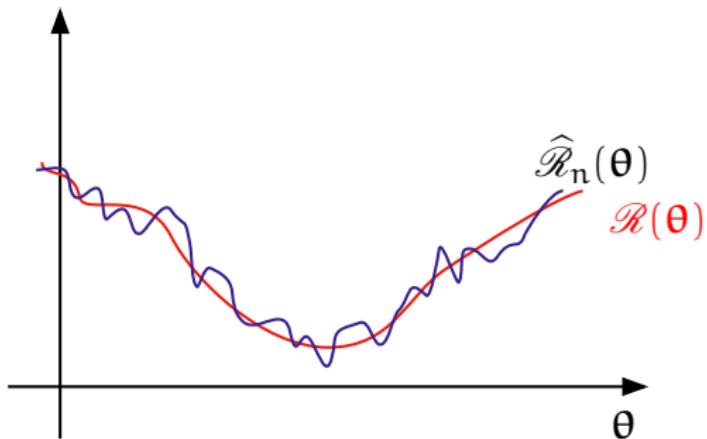
$$\max_{f \in \mathcal{F}_{\text{2lnn}}(\gamma)} |\widehat{\mathcal{R}}_n(f) - \mathcal{R}(f)| \lesssim \gamma \sqrt{\frac{d}{n}}. \quad \text{Independent of } N!$$

Statistical learning: Implications



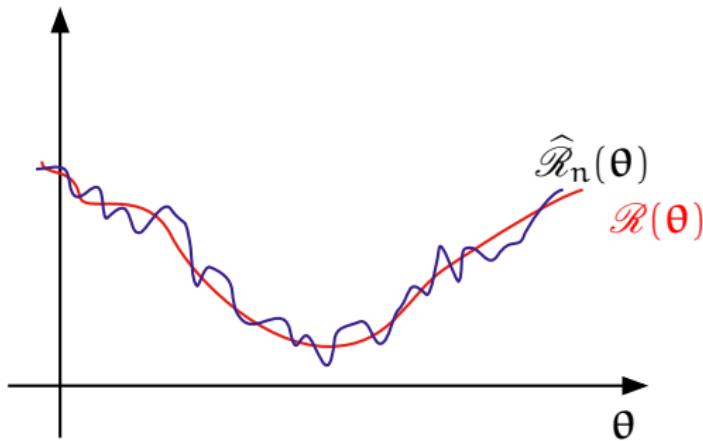
- ▶ Want to minimize $\mathcal{R}(\theta) := \mathbb{E}\{\ell(y_{n+1}, f(x_{n+1}; \theta))\}$ $\ell = \text{dist}$
- ▶ We have access to $\hat{\mathcal{R}}_n(\theta) := n^{-1} \sum_{i \leq n} \ell(y_i, f(x_i; \theta))$

Statistical learning: Implications



- Vapnik–Chervonenkis

Statistical learning: Implications

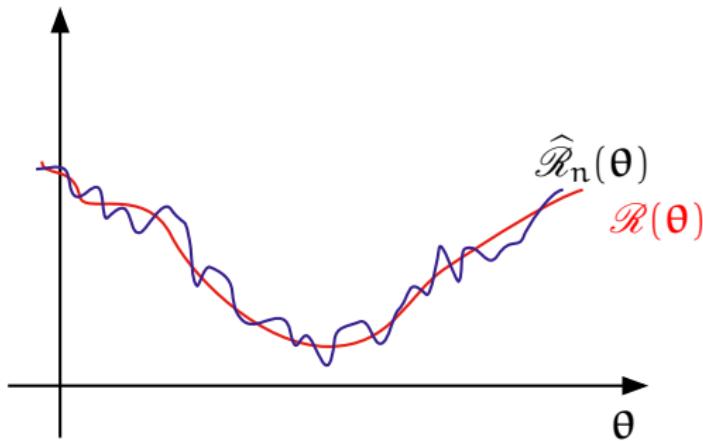


- ▶ Vapnik–Chervonenkis:

$$\max_{\theta \in \mathcal{S}} |\hat{\mathcal{R}}_n(\theta) - \mathcal{R}(\theta)| = O \left(\sqrt{\frac{d_{VC}(\mathcal{F})}{n}} \right)$$

- ▶ ⇒ Can minimize $\hat{\mathcal{R}}_n(\theta)$ instead of $\mathcal{R}(\theta)$
- ▶ Choose \mathcal{S} so that $d_{VC}(\mathcal{S}) \ll n$

Statistical learning: Implications

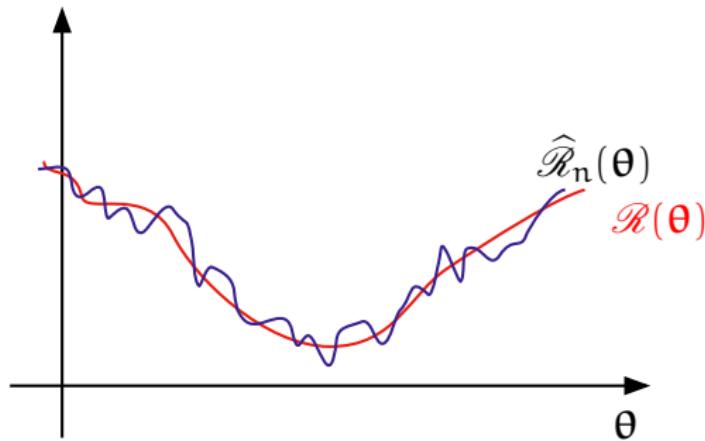


- ▶ Vapnik–Chervonenkis:

$$\max_{\theta \in \mathcal{S}} |\mathcal{R}(\theta) - \hat{\mathcal{R}}_n(\theta)| = O \left(\sqrt{\frac{d_{VC}(\mathcal{F})}{n}} \right)$$

- ▶ ⇒ Can minimize $\hat{\mathcal{R}}_n(\theta)$ instead of $\mathcal{R}(\theta)$
- ▶ Choose \mathcal{S} so that $d_{VC}(\mathcal{S}) \ll n$

Statistical learning: Implications

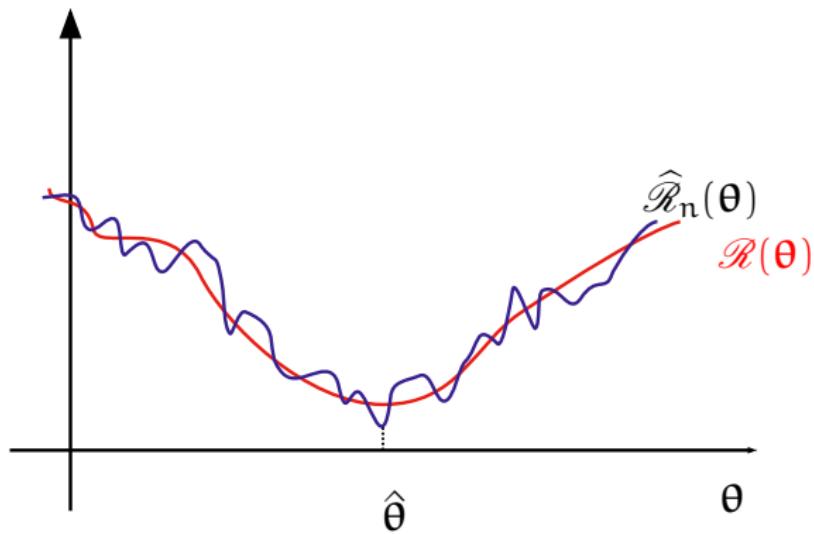


- ▶ Vapnik–Chervonenkis:

$$\max_{\theta \in \mathcal{S}} |\mathcal{R}(\theta) - \hat{\mathcal{R}}_n(\theta)| = O \left(\sqrt{\frac{d_{VC}(\mathcal{F})}{n}} \right)$$

- ▶ ⇒ Can minimize $\hat{\mathcal{R}}_n(\theta)$ instead of $\mathcal{R}(\theta)$
- ▶ Choose \mathcal{S} so that $d_{VC}(\mathcal{S}) \ll n$

Statistical learning

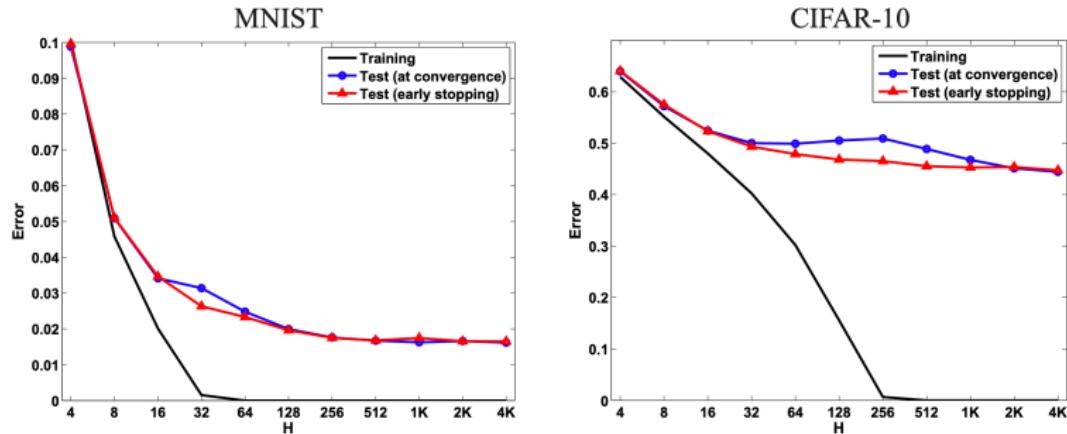


- Learned model $\hat{\theta}$

$$\text{Generalization Error} := \mathcal{R}(\hat{\theta}) - \hat{\mathcal{R}}_n(\hat{\theta})$$

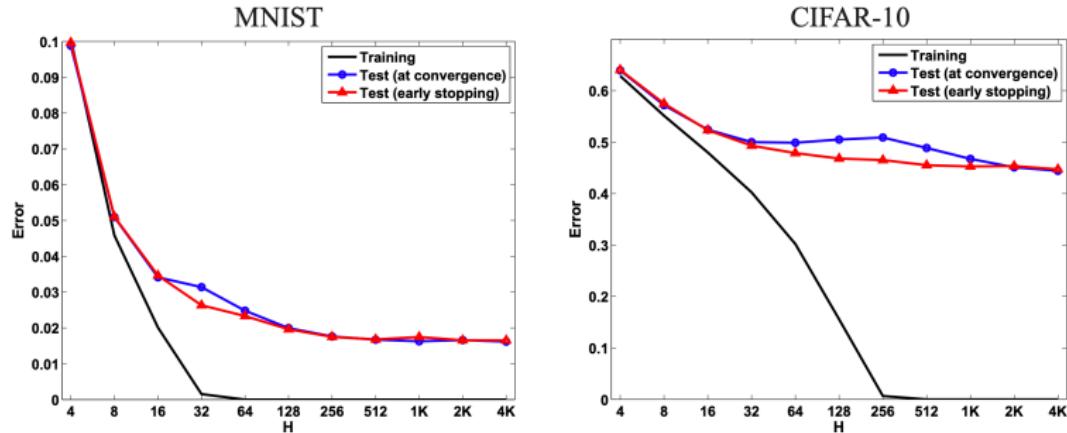
Enters modern machine learning

A small experiment



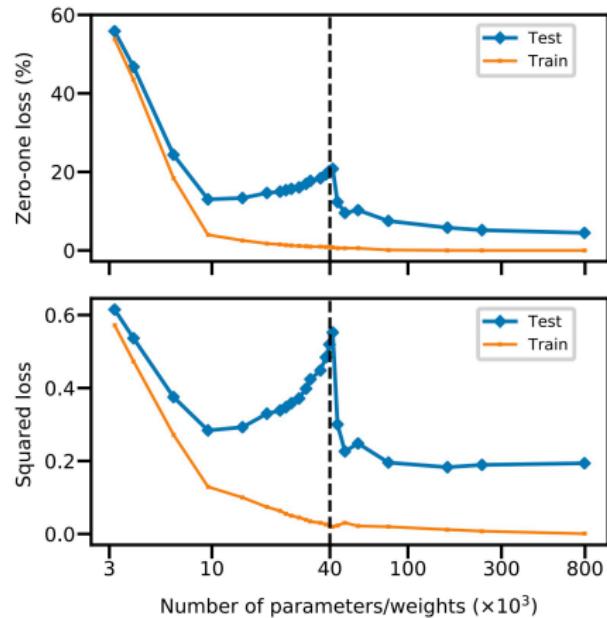
- ▶ x-axis \propto Number of parameters
- ▶ Converges to vanishing training error
- ▶ Resulting weights depend on the initialization

A small experiment



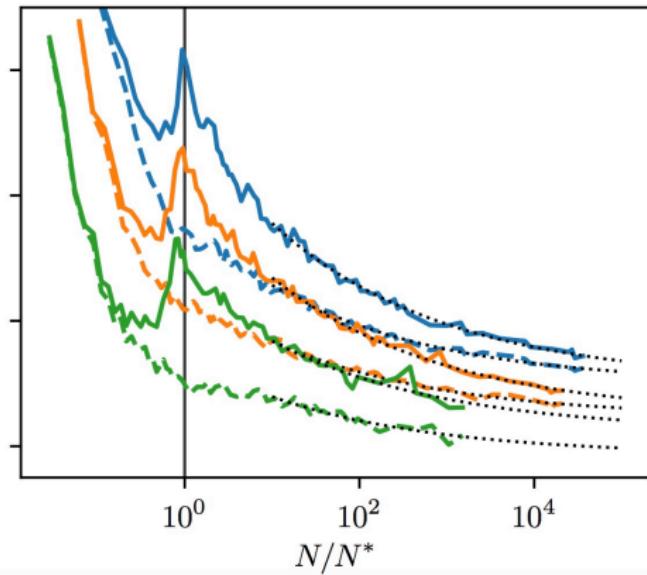
- ▶ x-axis \propto Number of parameters
- ▶ Converges to vanishing training error
- ▶ Resulting weights depend on the initialization

Another version of the same experiment



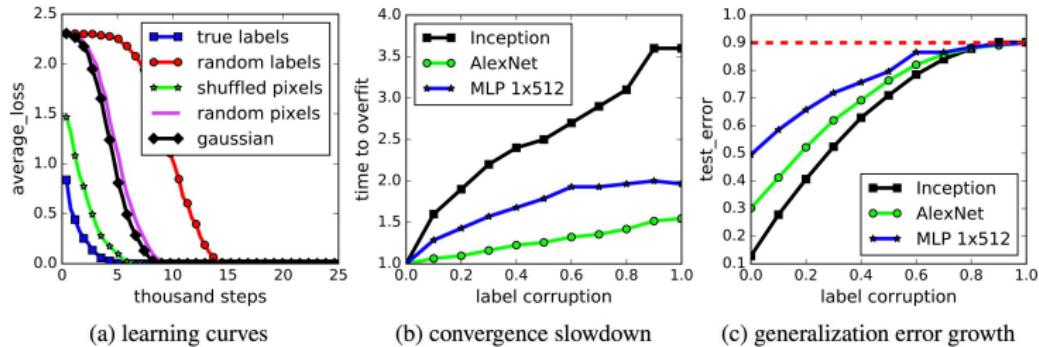
MNIST: 4,000 images 10 classes; 2-layers. Square loss. Belkin, Hsu, Ma, Mandal, 2018

Yet another version



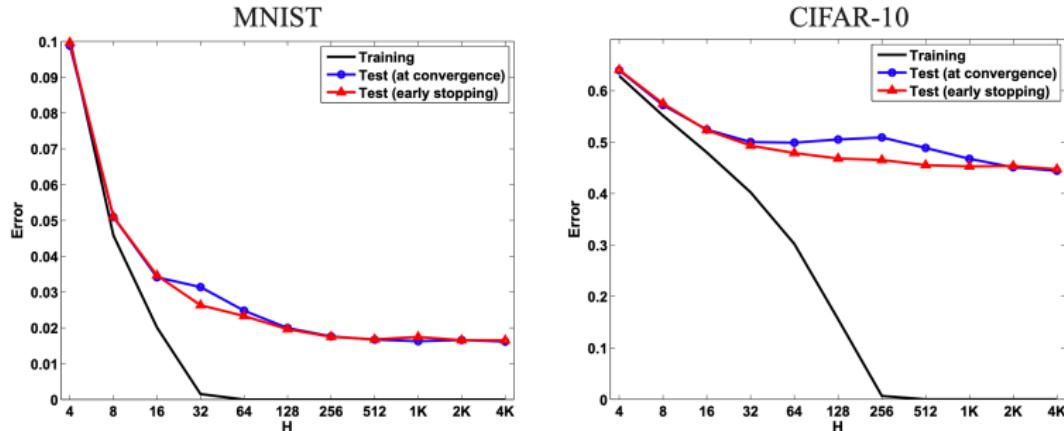
MNIST: 50,000 images in 2 different classes. 5-layers Neural Net. Quadratic hinge loss. Spigler,
Geiger, d'Ascoli, Sagun, Biroli, Wyart, 2018

A larger scale experiment



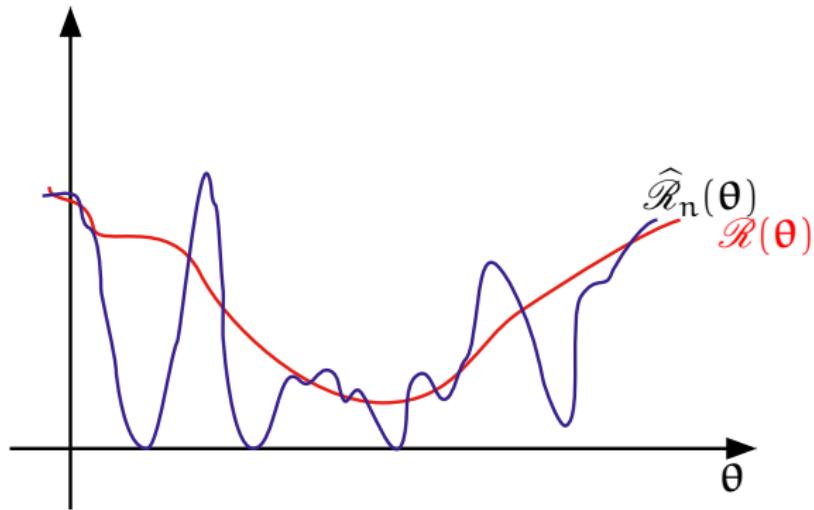
- ▶ Model complex enough to ‘interpolate’ random labels
- ▶ Despite this, does well on uncorrupted test samples
- ▶ Test error \gg Train error ≈ 0

What does this tell us about the landscape?



- ▶ Many (near-)global minima with $\hat{\mathcal{R}}_n(\boldsymbol{\theta}) \approx 0$.
- ▶ Close to a random initialization
- ▶ At global minima $0 \approx \hat{\mathcal{R}}_n(\boldsymbol{\theta}) \ll \mathcal{R}(\boldsymbol{\theta})$.

The actual landscape: A cartoon

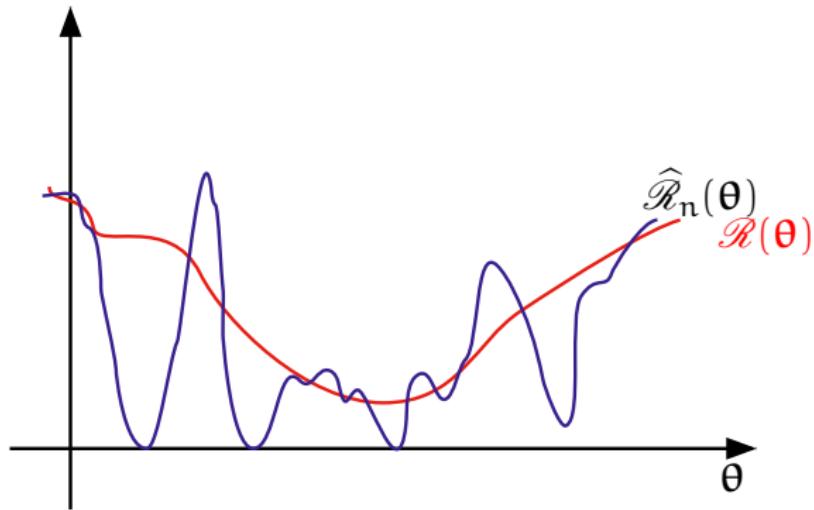


- ▶ How does generalization work?
- ▶ Why does minimizing $\hat{\mathcal{R}}_n(\theta)$ yields models with small $\mathcal{R}(\theta)$?

Many ‘bad’ global optima

Explanation must involve the optimization dynamics

The actual landscape: A cartoon

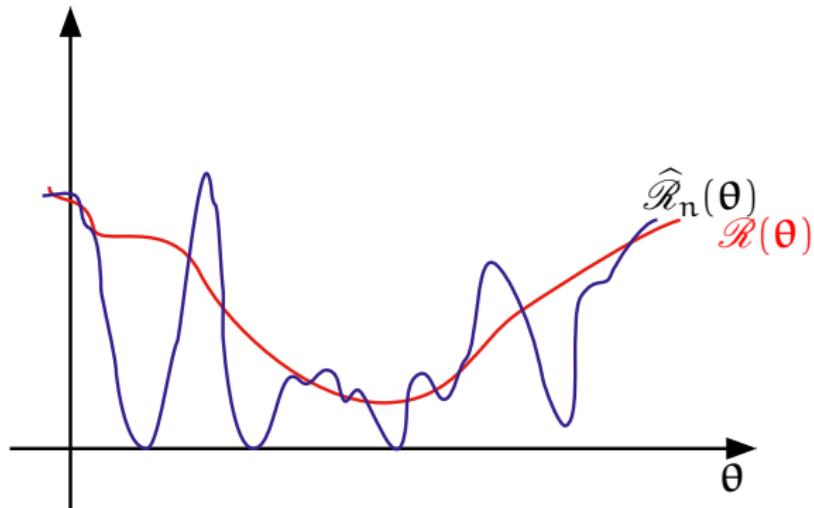


- ▶ How does generalization work?
- ▶ Why does minimizing $\hat{\mathcal{R}}_n(\theta)$ yields models with small $\mathcal{R}(\theta)$?

Many ‘bad’ global optima

Explanation must involve the optimization dynamics

The actual landscape: A cartoon

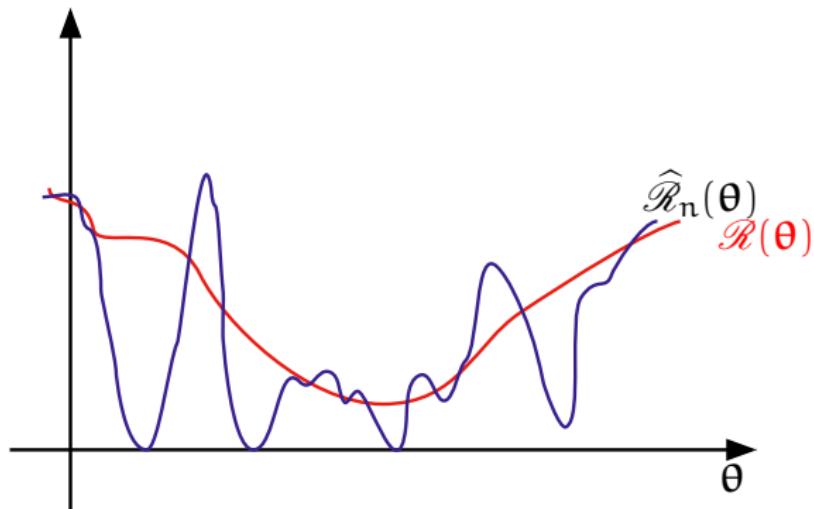


- ▶ How does generalization work?
- ▶ Why does minimizing $\hat{\mathcal{R}}_n(\theta)$ yields models with small $\mathcal{R}(\theta)$?

Many ‘bad’ global optima

Explanation must involve the optimization dynamics

The actual landscape



Two theories

Neural tangent theory:

Global minima near a random initialization behave well

Feature learning:

Few steps of gradient descent sufficient to converge close to $\operatorname{argmin}_{\theta} \mathcal{R}(\theta)$. (Overfitting cannot take place so quickly.)

Two-layer neural networks

$$f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle), \quad \boldsymbol{\theta} = ((a_i)_{i \leq N}, (\mathbf{w}_i)_{i \leq N}) \in \mathbb{R}^N \times (\mathbb{S}^{d-1})^N.$$

- ▶ $p = N d$
- ▶ Square loss train error:

$$\hat{\mathcal{R}}_n(\boldsymbol{\theta}) := \frac{1}{2n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2.$$

- ▶ Square loss test error:

$$\mathcal{R}(\boldsymbol{\theta}) := \frac{1}{2n} \mathbb{E}\{(y - f(\mathbf{x}; \boldsymbol{\theta}))^2\}.$$

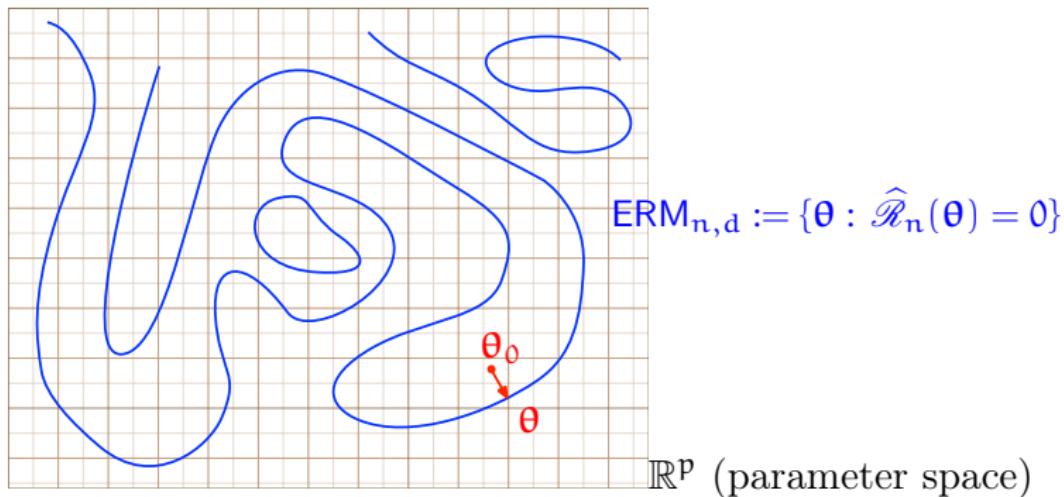
Single index model

(‘hydrogen atom’?)

Data $\{(\mathbf{x}_i, y_i) : i \leq n\}$ iid, $\varepsilon_i \sim N(0, \tau^2)$

$$\mathbf{x}_i \sim N(0, \mathbf{I}_d), \quad y_i = \varphi(\langle \mathbf{w}_*, \mathbf{x}_i \rangle) + \varepsilon_i$$

'Neural tangent' theory

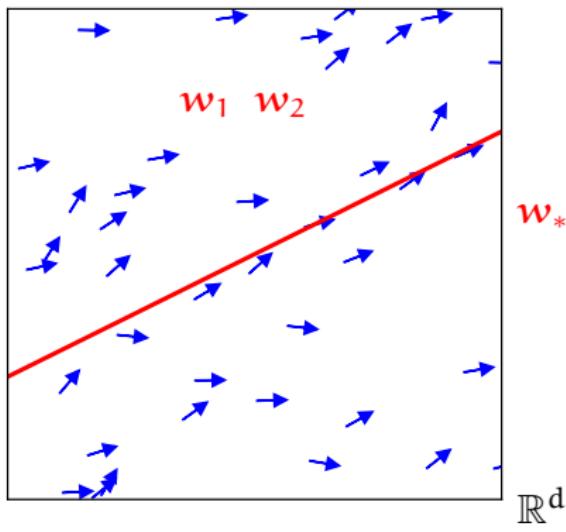


- ▶ Fast convergence to $\hat{\mathcal{R}}_n(\hat{\theta}) \approx 0$
- ▶ Overfitting $\hat{\mathcal{R}}_n(\hat{\theta}) \ll \mathcal{R}(\hat{\theta})$
- ▶ Fits min-norm kernel regression (see next)
- ▶ Statistically suboptimal.

Jacot, Gabriel, Hongler, 2018; Du, Zhai, Poczos, Singh 2018; Allen-Zhu, Li, Song 2018; Chizat, Bach, 2019; Oymak, Soltanolkotabi, 2019; ...

Feature learning

$$f(x; \theta) = \sum_{j=1}^N a_j \sigma(\langle w_j, x \rangle)$$



- ▶ Converges (ideally) to Bayes error $\hat{\mathcal{R}}_n(\hat{\theta}) \approx \tau^2/2$.
- ▶ No overfitting $\hat{\mathcal{R}}_n(\hat{\theta}) \approx \mathcal{R}(\hat{\theta})$
- ▶ Learns latent direction w_* : nonlinear.
- ▶ Expected to be statistically optimal.

Two theories

Neural tangent theory:

Global minima near a random initialization behave well

Feature learning:

Few steps of GD sufficient to converge to $\operatorname{argmin}_{\theta} \mathcal{R}(\theta)$.

Punchline: Each is correct in a different dynamical regime

Two theories

Neural tangent theory:

Global minima near a random initialization behave well

Feature learning:

Few steps of GD sufficient to converge to $\operatorname{argmin}_{\theta} \mathcal{R}(\theta)$.

Punchline: Each is correct in a different dynamical regime

Two previews

Preview #1: Generalization in the neural tangent regime

Neural tangent model for 2-layer networks

Linearize around initialization \mathbf{W}^0

$$f(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)$$

$$f_{NT}(\mathbf{x}; \bar{\mathbf{b}}, \mathbf{b}) = \sum_{j=1}^N \langle \mathbf{b}_j, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_j^0, \mathbf{x} \rangle) + \sum_{j=1}^N \bar{b}_j \sigma(\langle \mathbf{w}_j^0, \mathbf{x} \rangle).$$

Gradient descent \rightarrow closest interpolant

$$\hat{\mathbf{b}} = \operatorname{argmin} \left\{ \|\mathbf{b}\| : y_i = f_{NT}(\mathbf{x}_i; \mathbf{b}) \quad \forall i \leq n \right\}$$

Neural tangent model for 2-layer networks

Linearize around initialization \mathbf{W}^0

$$f(\mathbf{x}; \mathbf{a}, \mathbf{W}) = \sum_{j=1}^N a_j \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle)$$

$$f_{NT}(\mathbf{x}; \bar{\mathbf{b}}, \mathbf{b}) = \sum_{j=1}^N \langle \mathbf{b}_j, \mathbf{x} \rangle \sigma'(\langle \mathbf{w}_j^0, \mathbf{x} \rangle) + \sum_{j=1}^N \bar{b}_j \sigma(\langle \mathbf{w}_j^0, \mathbf{x} \rangle).$$

Gradient descent \longrightarrow closest interpolant

$$\hat{\mathbf{b}} = \operatorname{argmin} \left\{ \|\mathbf{b}\| : y_i = f_{NT}(\mathbf{x}_i; \mathbf{b}) \quad \forall i \leq n \right\}$$

Neural tangent model for 2-layer networks

Gradient descent \longrightarrow closest interpolant

$$\hat{\mathbf{b}} = \operatorname{argmin}_{\mathbf{b}} \left\{ \|\mathbf{b}\| : y_i = f_{\text{NT}}(x_i; \mathbf{b}) \quad \forall i \leq n \right\}$$

Equivalently

$$z_i = \Phi(x_i) := (x_i^\top \sigma'(\langle w_1^0, x_i \rangle), x_i^\top \sigma'(\langle w_2^0, x_i \rangle), \dots, x_i^\top \sigma'(\langle w_N^0, x_i \rangle))^T,$$

$$\hat{\mathbf{b}}_\lambda = \operatorname{argmin}_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_2^2 \right\}, \quad \lambda \downarrow 0.$$

Neural tangent model for 2-layer networks

Gradient descent \longrightarrow closest interpolant

$$\hat{\mathbf{b}} = \operatorname{argmin}_{\mathbf{b}} \left\{ \|\mathbf{b}\| : y_i = f_{\text{NT}}(x_i; \mathbf{b}) \quad \forall i \leq n \right\}$$

Equivalently

$$z_i = \Phi(x_i) := (x_i^T \sigma'(\langle w_1^0, x_i \rangle), x_i^T \sigma'(\langle w_2^0, x_i \rangle), \dots, x_i^T \sigma'(\langle w_N^0, x_i \rangle))^T,$$

$$\hat{\mathbf{b}}_\lambda = \operatorname{argmin}_{\mathbf{b}} \left\{ \|\mathbf{y} - \mathbf{Z}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_2^2 \right\}, \quad \lambda \downarrow 0.$$

Simpler setting: Random features model

$$\min_{\mathbf{b}} \left\{ \sum_{i=1}^n (y_i - f_{RF}(\mathbf{x}_i; \mathbf{b}))^2 + \lambda \|\mathbf{b}\|_2^2 \right\}, \quad f_{RF}(\mathbf{x}; \mathbf{b}) = \sum_{i=1}^N b_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle).$$

Theorem (Mei, M. 2019)

Assume $n, N, d \rightarrow \infty$

$$\frac{N}{d} \rightarrow \psi_1, \quad \frac{n}{d} \rightarrow \psi_2.$$

Decompose $\sigma(x) = \sigma_0 + \sigma_1 x + \sigma^{NL}(x)$ where (for $G \sim N(0, 1)$)

$\mathbb{E}[G\sigma^{NL}(G)] = \mathbb{E}[\sigma^{NL}(G)] = 0$, $\zeta^2 := \frac{\sigma_1^2}{\mathbb{E}[\sigma^{NL}(G)^2]}$. Then there are explicit functions $\mathcal{B}(\zeta, \psi_1, \psi_2, \bar{\lambda})$, $\mathcal{V}(\zeta, \psi_1, \psi_2, \bar{\lambda})$, such that

$$R(\hat{f}_\lambda) = F_1^2 \mathcal{B}(\zeta, \psi_1, \psi_2, \lambda) + (\tau^2 + F_*^2) \mathcal{V}(\zeta, \psi_1, \psi_2, \lambda) + F_*^2 + o_d(1).$$

Simpler setting: Random features model

$$\min_{\mathbf{b}} \left\{ \sum_{i=1}^n (y_i - f_{RF}(\mathbf{x}_i; \mathbf{b}))^2 + \lambda \|\mathbf{b}\|_2^2 \right\}, \quad f_{RF}(\mathbf{x}; \mathbf{b}) = \sum_{i=1}^N b_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle).$$

Theorem (Mei, M. 2019)

Assume $n, N, d \rightarrow \infty$

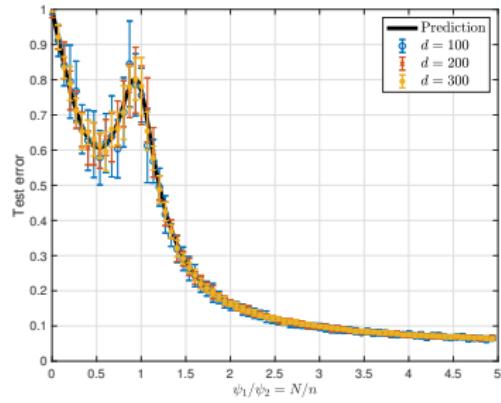
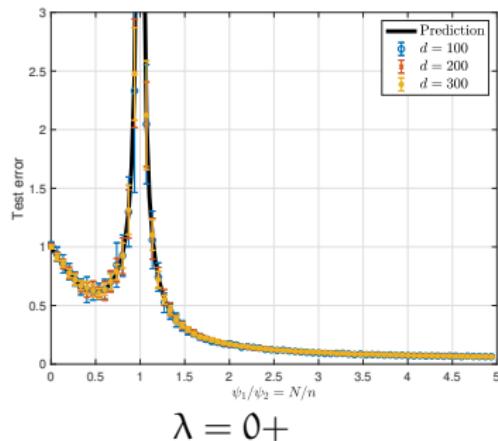
$$\frac{N}{d} \rightarrow \psi_1, \quad \frac{n}{d} \rightarrow \psi_2.$$

Decompose $\sigma(x) = \sigma_0 + \sigma_1 x + \sigma^{NL}(x)$ where (for $G \sim N(0, 1)$)

$\mathbb{E}[G\sigma^{NL}(G)] = \mathbb{E}[\sigma^{NL}(G)] = 0$, $\zeta^2 := \frac{\sigma_1^2}{\mathbb{E}[\sigma^{NL}(G)^2]}$. Then there are explicit functions $\mathcal{B}(\zeta, \psi_1, \psi_2, \bar{\lambda})$, $\mathcal{V}(\zeta, \psi_1, \psi_2, \bar{\lambda})$, such that

$$R(\hat{f}_\lambda) = F_1^2 \mathcal{B}(\zeta, \psi_1, \psi_2, \lambda) + (\tau^2 + F_*^2) \mathcal{V}(\zeta, \psi_1, \psi_2, \lambda) + F_*^2 + o_d(1).$$

Insight: $N/n \gg 1$ optimal



- ▶ Solid line: Theoretical prediction
- ▶ Mathematics tool: Random matrix theory

Preview #2: Dynamical decoupling of feature learning and overfitting

Gradient flow dynamics

Gradient flow

$$\dot{\theta}(t) = -\frac{n}{d} \nabla \hat{\mathcal{R}}_n(\theta(t)).$$

Square loss train error

$$\hat{\mathcal{R}}_n(\theta) := \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i; \theta))^2, \quad \theta = (\mathbf{a}, \mathbf{W}).$$

Two-layer network

$$f(x; \theta) = \frac{1}{N} \sum_{j=1}^N a_j \sigma(\langle w_j, x \rangle).$$

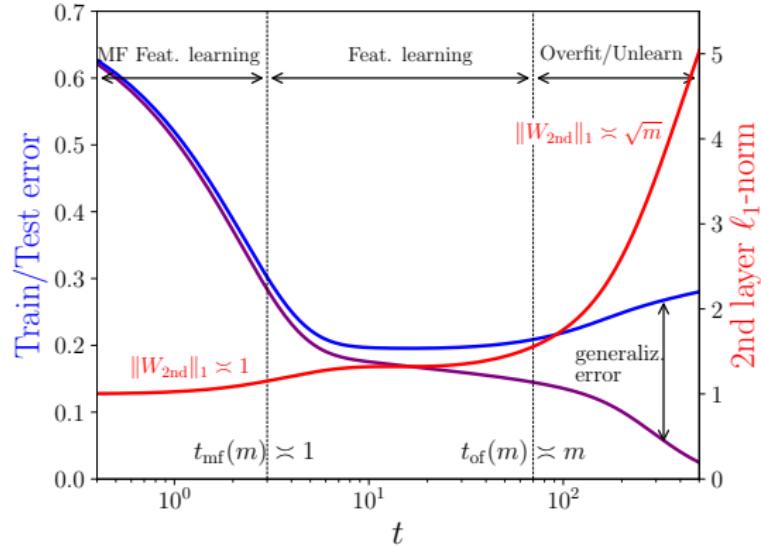
Dynamical approach

Gradient flow

$$\dot{\theta}(t) = -\frac{n}{d} \nabla \hat{\mathcal{R}}_n(\theta(t)).$$

1. Dynamics in a random landscape
2. Exact asymptotics as $n, d \rightarrow \infty$, $n/d \rightarrow \bar{\alpha}$ by rigorizing tools from theoretical physics
[DMFT, dynamical mean-field theory: Celentano, Cheng, M, 2022]
3. Large width limit $N, \bar{\alpha} \rightarrow \infty$, $\bar{\alpha}/N = \alpha$:
[M, Urbani, 2025]

Dynamical decoupling: Learning → Overfitting



$t \asymp 1$:

- ▶ Test error \approx Train error
- ▶ Learns latent direction \mathbf{w}_*
- ▶ Mean field asymptotically correct (Mei, M, Nguyen, 2018; Chizat, Bach, 2018; Rotskoff, Vanden Eijnden, 2018)

$t \asymp m \equiv N$:

- ▶ Test error \gg Train error; Train error $\rightarrow 0$
- ▶ *Unlearns* latent direction \mathbf{w}_*
- ▶ Neural tangent kernel learning

Conclusion

Conclusion

- ▶ Important to ask fundamental questions about AI.
- ▶ What does it mean that ‘AI models learn’?
- ▶ Requires to understand training dynamics

Thank you!

Conclusion

- ▶ Important to ask fundamental questions about AI.
- ▶ What does it mean that ‘AI models learn’?
- ▶ Requires to understand training dynamics

Thank you!