

Phase Transitions for Feature Learning in Neural Networks

Zihao Wang (Stanford University)

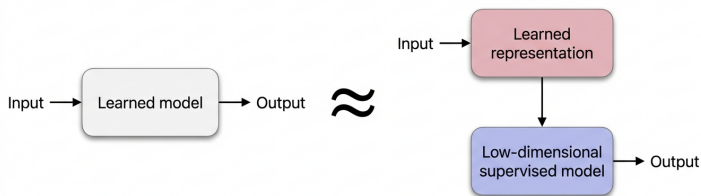
Center for Computational Mathematics, Flatiron Institute
March 2026



Joint work with Andrea Montanari
(Stanford University)

Folklore

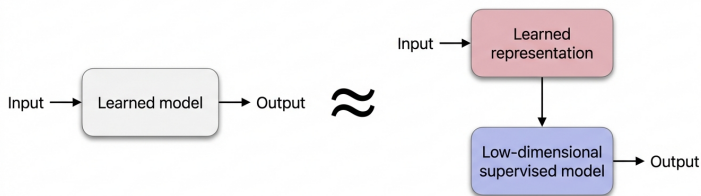
Huge success of deep learning from learning representations of data.



First learn representations, then fit a low-dimensional model over those representations!

Folklore

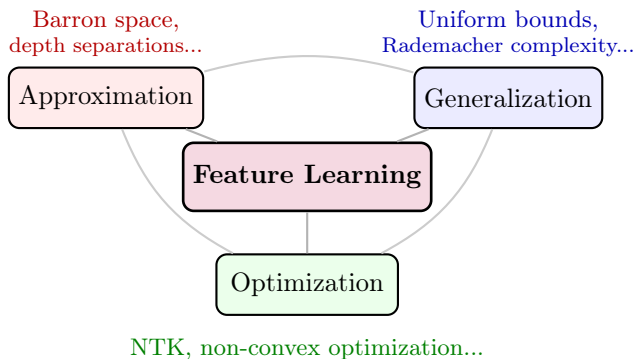
Huge success of deep learning from learning representations of data.



First learn representations, then fit a low-dimensional model over those representations!

Feature learning

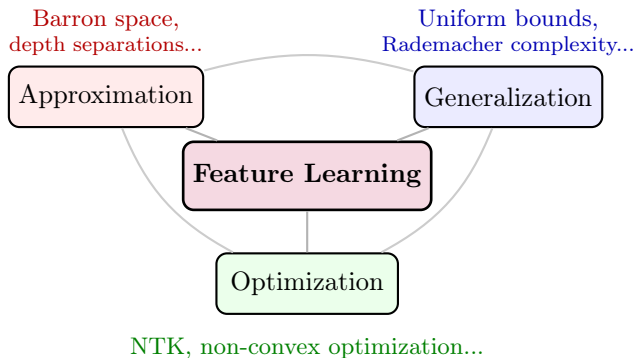
- ▶ Feature learning: Learning those representations from data.



In the feature learning regime, approximation, generalization, and optimization are **entangled** with each other.

Feature learning

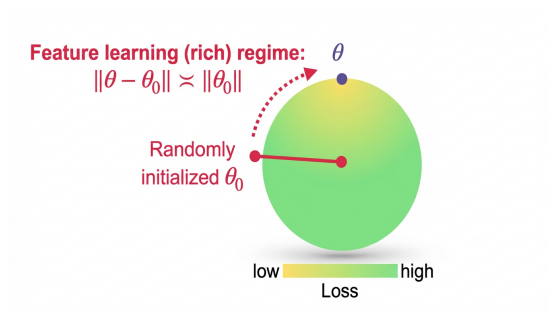
- ▶ Feature learning: Learning those representations from data.



In the feature learning regime, approximation, generalization, and optimization are **entangled** with each other.

Feature learning: Beyond NTK regime

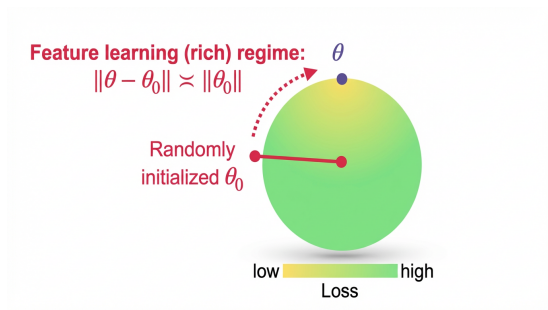
- ▶ Neural network $f : \mathcal{X} \times \mathbb{R}^D \rightarrow \mathbb{R}$.
Parameters $\theta \in \mathbb{R}^D$.



- ▶ Kernel/NTK/Lazy: $\|\theta - \theta_0\| \ll \|\theta_0\|$.

Feature learning: Beyond NTK regime

- ▶ Neural network $f : \mathcal{X} \times \mathbb{R}^D \rightarrow \mathbb{R}$.
Parameters $\theta \in \mathbb{R}^D$.



- ▶ Kernel/NTK/Lazy: $\|\theta - \theta_0\| \ll \|\theta_0\|$.

Research questions

- ▶ (Optimization) How models learn features through gradient descent? Understand **loss landscape and training dynamics**.
- ▶ (Generalization) The required **sample size** to learn features. How does this scale with target task, loss function, activation, etc.?

Research questions

- ▶ (Optimization) How models learn features through gradient descent? Understand **loss landscape and training dynamics**.
- ▶ (Generalization) The required **sample size** to learn features. How does this scale with target task, loss function, activation, etc.?

Outline

- 1 Setting and results overview
- 2 Numerical experiments
- 3 Formal results and outline of the proofs
- 4 Discussion and future directions

Multi-index models

Consider multi-index models as our targets for feature learning.

Setup: Observe n i.i.d. samples $(\mathbf{x}_i, y_i)_{i \leq n}$:

- ▶ Covariates: $\mathbf{x}_i \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$
- ▶ Responses:

$$y_i = h(\Theta_*^\top \mathbf{x}_i, \varepsilon_i), \quad \varepsilon_i \sim \mathbf{N}(0, 1)$$

- ▶ $\Theta_* \in \mathbb{R}^{d \times k}$ ($k \ll d$)
- ▶ $h: \mathbb{R}^{k+1} \rightarrow \mathbb{R}$ is the link function.

Multi-index models

Consider multi-index models as our targets for feature learning.

Setup: Observe n i.i.d. samples $(\mathbf{x}_i, y_i)_{i \leq n}$:

- ▶ Covariates: $\mathbf{x}_i \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$
- ▶ Responses:

$$y_i = h(\Theta_*^\top \mathbf{x}_i, \varepsilon_i), \quad \varepsilon_i \sim \mathbf{N}(0, 1)$$

- ▶ $\Theta_* \in \mathbb{R}^{d \times k}$ ($k \ll d$)
- ▶ $h: \mathbb{R}^{k+1} \rightarrow \mathbb{R}$ is the link function.

Multi-index models

Consider multi-index models as our targets for feature learning.

Setup: Observe n i.i.d. samples $(\mathbf{x}_i, y_i)_{i \leq n}$:

- ▶ Covariates: $\mathbf{x}_i \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$
- ▶ Responses:

$$y_i = h(\Theta_*^\top \mathbf{x}_i, \varepsilon_i), \quad \varepsilon_i \sim \mathbf{N}(0, 1)$$

- ▶ $\Theta_* \in \mathbb{R}^{d \times k}$ ($k \ll d$)
- ▶ $h: \mathbb{R}^{k+1} \rightarrow \mathbb{R}$ is the link function.

Multi-index models

Consider multi-index models as our targets for feature learning.

Setup: Observe n i.i.d. samples $(\mathbf{x}_i, y_i)_{i \leq n}$:

- ▶ Covariates: $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$
- ▶ Responses:

$$y_i = h(\Theta_*^\top \mathbf{x}_i, \varepsilon_i), \quad \varepsilon_i \sim \mathcal{N}(0, 1)$$

- ▶ $\Theta_* \in \mathbb{R}^{d \times k}$ ($k \ll d$)
- ▶ $h : \mathbb{R}^{k+1} \rightarrow \mathbb{R}$ is the link function.

Multi-index models

Consider multi-index models as our targets for feature learning.

Setup: Observe n i.i.d. samples $(\mathbf{x}_i, y_i)_{i \leq n}$:

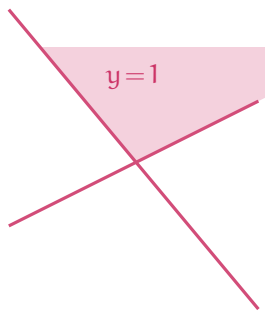
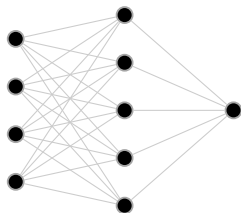
- ▶ Covariates: $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$
- ▶ Responses:

$$y_i = h(\Theta_*^\top \mathbf{x}_i, \varepsilon_i), \quad \varepsilon_i \sim \mathcal{N}(0, 1)$$

- ▶ $\Theta_* \in \mathbb{R}^{d \times k}$ ($k \ll d$)
- ▶ $h : \mathbb{R}^{k+1} \rightarrow \mathbb{R}$ is the link function.

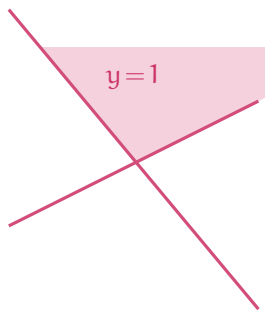
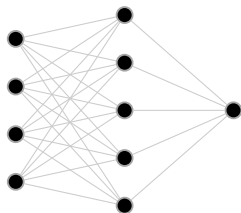
Examples of multi-index models

- ▶ **Single-index** ($k = 1$):
(e.g., phase retrieval) $y = (\boldsymbol{\theta}_*^T \mathbf{x})^2$
- ▶ **Neural networks** ($O(1)$ neurons):
 $y = \sum_{j=1}^k a_j \sigma(\boldsymbol{\theta}_{*j}^T \mathbf{x}) + \varepsilon$
- ▶ **Intersection of halfspaces**:
 $y = \prod_{j=1}^k \mathbf{1}_{\{\boldsymbol{\theta}_{*j}^T \mathbf{x} > 0\}}$
- ▶ **Polynomials on $O(1)$ directions**:
 $y = p(\boldsymbol{\Theta}_*^T \mathbf{x})$



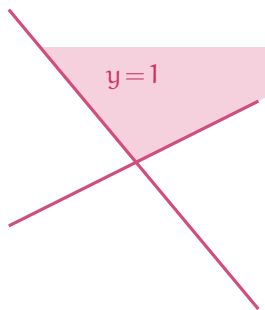
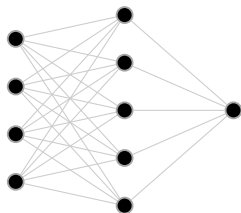
Examples of multi-index models

- ▶ **Single-index** ($k = 1$):
(e.g., phase retrieval) $y = (\boldsymbol{\theta}_*^\top \mathbf{x})^2$
- ▶ **Neural networks** ($O(1)$ neurons):
 $y = \sum_{j=1}^k \alpha_j \sigma(\boldsymbol{\theta}_{*j}^\top \mathbf{x}) + \varepsilon$
- ▶ **Intersection of halfspaces**:
 $y = \prod_{j=1}^k \mathbf{1}_{\{\boldsymbol{\theta}_{*j}^\top \mathbf{x} > 0\}}$
- ▶ **Polynomials on $O(1)$ directions**:
 $y = p(\boldsymbol{\Theta}_*^\top \mathbf{x})$



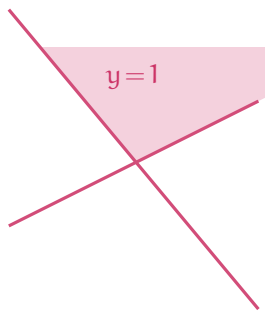
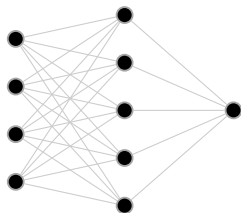
Examples of multi-index models

- ▶ **Single-index** ($k = 1$):
(e.g., phase retrieval) $y = (\boldsymbol{\theta}_*^\top \mathbf{x})^2$
- ▶ **Neural networks** ($O(1)$ neurons):
 $y = \sum_{j=1}^k \alpha_j \sigma(\boldsymbol{\theta}_{*j}^\top \mathbf{x}) + \varepsilon$
- ▶ **Intersection of halfspaces**:
 $y = \prod_{j=1}^k \mathbf{1}_{\{\boldsymbol{\theta}_{*j}^\top \mathbf{x} > 0\}}$
- ▶ **Polynomials on $O(1)$ directions**:
 $y = p(\boldsymbol{\Theta}_*^\top \mathbf{x})$



Examples of multi-index models

- ▶ **Single-index** ($k = 1$):
(e.g., phase retrieval) $y = (\boldsymbol{\theta}_*^\top \mathbf{x})^2$
- ▶ **Neural networks** ($O(1)$ neurons):
 $y = \sum_{j=1}^k \alpha_j \sigma(\boldsymbol{\theta}_{*j}^\top \mathbf{x}) + \varepsilon$
- ▶ **Intersection of halfspaces**:
 $y = \prod_{j=1}^k \mathbf{1}_{\{\boldsymbol{\theta}_{*j}^\top \mathbf{x} > 0\}}$
- ▶ **Polynomials** on $O(1)$ directions:
 $y = p(\boldsymbol{\Theta}_*^\top \mathbf{x})$



The learning problem

Feature learning = Learning the latent span(Θ_*)

Question: How many samples do we need to perform feature learning?

- ▶ Simple enough to be tractable mathematically
- ▶ Complex enough that NTK/RF/KRR cannot learn¹

Information-theoretic answer:² $n = \Theta(d)$

¹[Yehudai and Shamir '19], [Ghorbani, Mei, Misiakiewicz and Montanari '19]

²[Barbier, Krzakala, Macris+ '19], [Aubin, Maillard, Barbier+ '19]

The learning problem

Feature learning = Learning the latent span(Θ_*)

Question: How many samples do we need to perform feature learning?

- ▶ Simple enough to be tractable mathematically
- ▶ Complex enough that NTK/RF/KRR cannot learn¹

Information-theoretic answer:² $n = \Theta(d)$

¹[Yehudai and Shamir '19], [Ghorbani, Mei, Misiakiewicz and Montanari '19]

²[Barbier, Krzakala, Macris+ '19], [Aubin, Maillard, Barbier+ '19]

The learning problem

Feature learning = Learning the latent span(Θ_*)

Question: How many samples do we need to perform feature learning?

- ▶ Simple enough to be tractable mathematically
- ▶ Complex enough that NTK/RF/KRR cannot learn¹

Information-theoretic answer:² $n = \Theta(d)$

¹[Yehudai and Shamir '19], [Ghorbani, Mei, Misiakiewicz and Montanari '19]

²[Barbier, Krzakala, Macris+ '19], [Aubin, Maillard, Barbier+ '19]

The learning problem

Feature learning = Learning the latent span(Θ_*)

Question: How many samples do we need to perform feature learning?

- ▶ Simple enough to be tractable mathematically
- ▶ Complex enough that NTK/RF/KRR cannot learn¹

Information-theoretic answer:² $n = \Theta(d)$

¹[Yehudai and Shamir '19], [Ghorbani, Mei, Misiakiewicz and Montanari '19]

²[Barbier, Krzakala, Macris+ '19], [Aubin, Maillard, Barbier+ '19]

Prior work: Learning multi-index models

Insight: $n, d \rightarrow \infty, n/d \rightarrow \delta$.

- ▶ **Statistical threshold** δ_{IT} :
Below δ_{IT} , learning is information-theoretically impossible.
- ▶ **Algorithmic threshold** δ_{alg} :
Below δ_{alg} , no polynomial-time algorithm can learn Θ_* .

Gap: Sometimes $\delta_{\text{IT}} < \delta_{\text{alg}}$ — computational hardness

[Mondelli, Montanari '18], [Lu, Li '19], [Chen, Meka '20], [Damian, Pillaud-Vivien, Lee, Bruna '24], [Damian, Lee, Bruna '25], [Troiani, Dandi+ '25], [Kovačević, Zhang, Mondelli '25], [Defilippis, Dandi+ '25], ...

Neural Networks? (we want an analogous δ_{NN})

Prior work: Learning multi-index models

Insight: $n, d \rightarrow \infty, n/d \rightarrow \delta$.

- ▶ **Statistical threshold δ_{IT} :**
Below δ_{IT} , learning is information-theoretically impossible.
- ▶ **Algorithmic threshold δ_{alg} :**
Below δ_{alg} , no polynomial-time algorithm can learn Θ_* .

Gap: Sometimes $\delta_{\text{IT}} < \delta_{\text{alg}}$ — computational hardness

[Mondelli, Montanari '18], [Lu, Li '19], [Chen, Meka '20], [Damian, Pillaud-Vivien, Lee, Bruna '24], [Damian, Lee, Bruna '25], [Troiani, Dandi+ '25], [Kovačević, Zhang, Mondelli '25], [Defilippis, Dandi+ '25], ...

Neural Networks? (we want an analogous δ_{NN})

Prior work: Learning multi-index models

Insight: $n, d \rightarrow \infty, n/d \rightarrow \delta$.

- ▶ **Statistical threshold δ_{IT} :**
Below δ_{IT} , learning is information-theoretically impossible.
- ▶ **Algorithmic threshold δ_{alg} :**
Below δ_{alg} , no polynomial-time algorithm can learn Θ_* .

Gap: Sometimes $\delta_{\text{IT}} < \delta_{\text{alg}}$ — computational hardness

[Mondelli, Montanari '18], [Lu, Li '19], [Chen, Meka '20], [Damian, Pillaud-Vivien, Lee, Bruna '24], [Damian, Lee, Bruna '25], [Troiani, Dandi+ '25], [Kovačević, Zhang, Mondelli '25], [Defilippis, Dandi+ '25], ...

Neural Networks? (we want an analogous δ_{NN})

Prior work: Learning multi-index models

Insight: $n, d \rightarrow \infty, n/d \rightarrow \delta$.

- ▶ **Statistical threshold δ_{IT} :**
Below δ_{IT} , learning is information-theoretically impossible.
- ▶ **Algorithmic threshold δ_{alg} :**
Below δ_{alg} , no polynomial-time algorithm can learn Θ_* .

Gap: Sometimes $\delta_{\text{IT}} < \delta_{\text{alg}}$ — computational hardness

[Mondelli, Montanari '18], [Lu, Li '19], [Chen, Meka '20], [Damian, Pillaud-Vivien, Lee, Bruna '24], [Damian, Lee, Bruna '25], [Troiani, Dandi+ '25], [Kovačević, Zhang, Mondelli '25], [Defilippis, Dandi+ '25], ...

Neural Networks? (we want an analogous δ_{NN})

Prior work: Learning multi-index models

Insight: $n, d \rightarrow \infty, n/d \rightarrow \delta$.

- ▶ **Statistical threshold δ_{IT} :**
Below δ_{IT} , learning is information-theoretically impossible.
- ▶ **Algorithmic threshold δ_{alg} :**
Below δ_{alg} , no polynomial-time algorithm can learn Θ_* .

Gap: Sometimes $\delta_{\text{IT}} < \delta_{\text{alg}}$ — computational hardness

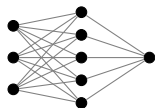
[Mondelli, Montanari '18], [Lu, Li '19], [Chen, Meka '20], [Damian, Pillaud-Vivien, Lee, Bruna '24], [Damian, Lee, Bruna '25], [Troiani, Dandi+ '25], [Kovačević, Zhang, Mondelli '25], [Defilippis, Dandi+ '25], ...

Neural Networks? (we want an analogous δ_{NN})

Neural network setup

Two-layer neural network:

$$f_{\Theta}(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m \mathbf{a}_j \sigma(\boldsymbol{\theta}_j^{\top} \mathbf{x} + \mathbf{b}_j)$$



- ▶ $\Theta = [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m] \in \mathbb{R}^{d \times m}$: first-layer weights
- ▶ σ : activation function (e.g., GeLU, ReLU, ...)
- ▶ $(\mathbf{a}_j, \mathbf{b}_j)$: fixed second-layer weights and biases

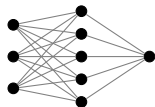
Training: Only train Θ via gradient descent on empirical risk

$$\text{Risk}(\Theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\Theta}(\mathbf{x}_i))$$

Neural network setup

Two-layer neural network:

$$f_{\Theta}(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m \mathbf{a}_j \sigma(\boldsymbol{\theta}_j^{\top} \mathbf{x} + \mathbf{b}_j)$$



- ▶ $\Theta = [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m] \in \mathbb{R}^{d \times m}$: first-layer weights
- ▶ σ : activation function (e.g., GeLU, ReLU, ...)
- ▶ $(\mathbf{a}_j, \mathbf{b}_j)$: fixed second-layer weights and biases

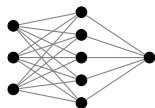
Training: Only train Θ via gradient descent on empirical risk

$$\text{Risk}(\Theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\Theta}(\mathbf{x}_i))$$

Neural network setup

Two-layer neural network:

$$f_{\Theta}(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\boldsymbol{\theta}_j^{\top} \mathbf{x} + b_j)$$



- ▶ $\Theta = [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m] \in \mathbb{R}^{d \times m}$: first-layer weights
- ▶ σ : activation function (e.g., GeLU, ReLU, ...)
- ▶ (a_j, b_j) : fixed second-layer weights and biases

Training: Only train Θ via gradient descent on empirical risk

$$\text{Risk}(\Theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\Theta}(\mathbf{x}_i))$$

Training protocol

Full-batch gradient descent:

$$\Theta(t+1) = \Theta(t) - \eta \nabla_{\Theta} \text{Risk}(\Theta(t))$$

Initialization: $\theta_j \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathbb{S}^{d-1})$

Proportional asymptotics:

$$n, d \rightarrow \infty, \quad n/d \rightarrow \delta \in (0, \infty)$$

Training protocol

Full-batch gradient descent:

$$\Theta(t+1) = \Theta(t) - \eta \nabla_{\Theta} \text{Risk}(\Theta(t))$$

Initialization: $\theta_j \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathbb{S}^{d-1})$

Proportional asymptotics:

$$n, d \rightarrow \infty, \quad n/d \rightarrow \delta \in (0, \infty)$$

Training protocol

Full-batch gradient descent:

$$\Theta(t+1) = \Theta(t) - \eta \nabla_{\Theta} \text{Risk}(\Theta(t))$$

Initialization: $\theta_j \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\mathbb{S}^{d-1})$

Proportional asymptotics:

$$n, d \rightarrow \infty, \quad n/d \rightarrow \delta \in (0, \infty)$$

Prior work: Feature learning in neural networks

[Ben Arous+ '21], [Damian+ '22], [Ba+ '22], [Abbe+ '23], [Bietti+ '23], [**Wang**+ '23], [Berthier+ '24], [Dandi+ '24], [Lee+ '24], [Fu, **Wang**+ '24], [Montanari, Urbani '25], [Zhang*, **Wang***+ '25], ...

- (*) [Ben Arous+ '21]: introduced 'information exponent'
- (*) [Abbe+ '23]: introduced 'leap complexity'
- (*) [Zhang*, **Wang***+ '25]: For generic³ multi-index models, $n = \tilde{O}(d)$ samples suffice for gradient descent on neural networks

³Generative leap exponent ≤ 2 with no generative-staircase structure; covers almost all common examples.

Prior work: Feature learning in neural networks

[Ben Arous+ '21], [Damian+ '22], [Ba+ '22], [Abbe+ '23], [Bietti+ '23], [**Wang**+ '23], [Berthier+ '24], [Dandi+ '24], [Lee+ '24], [Fu, **Wang**+ '24], [Montanari, Urbani '25], [Zhang*, **Wang***+ '25], ...

- (*) [Ben Arous+ '21]: introduced 'information exponent'
- (*) [Abbe+ '23]: introduced 'leap complexity'
- (*) [Zhang*, **Wang***+ '25]: For generic³ multi-index models, $n = \tilde{O}(d)$ samples suffice for gradient descent on neural networks

³Generative leap exponent ≤ 2 with no generative-staircase structure; covers almost all common examples.

Neural networks?

Question: What is the **learning threshold** δ_{NN} for gradient descent on neural networks?

- ▶ Do NNs achieve the same threshold as δ_{alg} ?
- ▶ If not, what is the gap? How does it depend on:
 - ▶ Architecture (activation, width, ...)
 - ▶ Algorithm (loss, initialization, stepsize, ...)
- ▶ What is the mechanism that determines δ_{NN} ?

Neural networks?

Question: What is the **learning threshold** δ_{NN} for gradient descent on neural networks?

- ▶ Do NNs achieve the same threshold as δ_{alg} ?
- ▶ If not, what is the gap? How does it depend on:
 - ▶ Architecture (activation, width, ...)
 - ▶ Algorithm (loss, initialization, stepsize, ...)
- ▶ What is the mechanism that determines δ_{NN} ?

Neural networks?

Question: What is the **learning threshold** δ_{NN} for gradient descent on neural networks?

- ▶ Do NNs achieve the same threshold as δ_{alg} ?
- ▶ If not, what is the gap? How does it depend on:
 - ▶ Architecture (activation, width, ...)
 - ▶ Algorithm (loss, initialization, stepsize, ...)
- ▶ What is the mechanism that determines δ_{NN} ?

Neural networks?

Question: What is the **learning threshold** δ_{NN} for gradient descent on neural networks?

- ▶ Do NNs achieve the same threshold as δ_{alg} ?
- ▶ If not, what is the gap? How does it depend on:
 - ▶ Architecture (activation, width, ...)
 - ▶ Algorithm (loss, initialization, stepsize, ...)
- ▶ What is the mechanism that determines δ_{NN} ?

Our results (informal overview)

Main contribution: Assuming the network width m is one or a large enough constant, given the knowledge of

- ▶ Target function $h(\cdot)$,
- ▶ Loss function $\ell(\cdot, \cdot)$,
- ▶ Activation function $\sigma(\cdot)$,
- ▶ Learning rate η

We derive an **exact, explicit formula** for computing δ_{NN} :

- ▶ $\delta > \delta_{\text{NN}}$: Hessian of empirical risk after $O(1)$ GD steps has informative⁴ low-lying eigenvectors
- ▶ $\delta < \delta_{\text{NN}}$: no informative eigenvector of the Hessian

δ_{NN} locates a **spectral phase transition**.

⁴Non-vanishing correlation with the ‘non-trivial’ directions of Θ_{**} .

Our results (informal overview)

Main contribution: Assuming the network width m is one or a large enough constant, given the knowledge of

- ▶ Target function $h(\cdot)$,
- ▶ Loss function $\ell(\cdot, \cdot)$,
- ▶ Activation function $\sigma(\cdot)$,
- ▶ Learning rate η

We derive an **exact, explicit formula** for computing δ_{NN} :

- ▶ $\delta > \delta_{\text{NN}}$: Hessian of empirical risk after $O(1)$ GD steps has informative⁴ low-lying eigenvectors
- ▶ $\delta < \delta_{\text{NN}}$: no informative eigenvector of the Hessian

δ_{NN} locates a **spectral phase transition**.

⁴Non-vanishing correlation with the ‘non-trivial’ directions of Θ_{**} .

Our results (informal overview)

Main contribution: Assuming the network width m is one or a large enough constant, given the knowledge of

- ▶ Target function $h(\cdot)$,
- ▶ Loss function $\ell(\cdot, \cdot)$,
- ▶ Activation function $\sigma(\cdot)$,
- ▶ Learning rate η

We derive an **exact, explicit formula** for computing δ_{NN} :

- ▶ $\delta > \delta_{\text{NN}}$: Hessian of empirical risk after $O(1)$ GD steps has informative⁴ low-lying eigenvectors
- ▶ $\delta < \delta_{\text{NN}}$: no informative eigenvector of the Hessian

δ_{NN} locates a **spectral phase transition**.

⁴Non-vanishing correlation with the ‘non-trivial’ directions of Θ_* .

Our results (informal overview)

Main contribution: Assuming the network width m is one or a large enough constant, given the knowledge of

- ▶ Target function $h(\cdot)$,
- ▶ Loss function $\ell(\cdot, \cdot)$,
- ▶ Activation function $\sigma(\cdot)$,
- ▶ Learning rate η

We derive an **exact, explicit formula** for computing δ_{NN} :

- ▶ $\delta > \delta_{\text{NN}}$: Hessian of empirical risk after $O(1)$ GD steps has informative⁴ low-lying eigenvectors
- ▶ $\delta < \delta_{\text{NN}}$: no informative eigenvector of the Hessian

δ_{NN} locates a **spectral phase transition**.

⁴Non-vanishing correlation with the ‘non-trivial’ directions of Θ_* .

Our results (informal overview)

Main contribution: Assuming the network width m is one or a large enough constant, given the knowledge of

- ▶ Target function $h(\cdot)$,
- ▶ Loss function $\ell(\cdot, \cdot)$,
- ▶ Activation function $\sigma(\cdot)$,
- ▶ Learning rate η

We derive an **exact, explicit formula** for computing δ_{NN} :

- ▶ $\delta > \delta_{\text{NN}}$: Hessian of empirical risk after $O(1)$ GD steps has informative⁴ low-lying eigenvectors
- ▶ $\delta < \delta_{\text{NN}}$: no informative eigenvector of the Hessian

δ_{NN} locates a **spectral phase transition**.

⁴Non-vanishing correlation with the ‘non-trivial’ directions of Θ_* .

Our results (informal overview)

Main contribution: Assuming the network width m is one or a large enough constant, given the knowledge of

- ▶ Target function $h(\cdot)$,
- ▶ Loss function $\ell(\cdot, \cdot)$,
- ▶ Activation function $\sigma(\cdot)$,
- ▶ Learning rate η

We derive an **exact, explicit formula** for computing δ_{NN} :

- ▶ $\delta > \delta_{\text{NN}}$: Hessian of empirical risk after $O(1)$ GD steps has informative⁴ low-lying eigenvectors
- ▶ $\delta < \delta_{\text{NN}}$: no informative eigenvector of the Hessian

δ_{NN} locates a **spectral phase transition**.

⁴Non-vanishing correlation with the ‘non-trivial’ directions of Θ_{**} .

Our results (informal overview)

Main contribution: Assuming the network width m is one or a large enough constant, given the knowledge of

- ▶ Target function $h(\cdot)$,
- ▶ Loss function $\ell(\cdot, \cdot)$,
- ▶ Activation function $\sigma(\cdot)$,
- ▶ Learning rate η

We derive an **exact, explicit formula** for computing δ_{NN} :

- ▶ $\delta > \delta_{\text{NN}}$: Hessian of empirical risk after $O(1)$ GD steps has informative⁴ low-lying eigenvectors
- ▶ $\delta < \delta_{\text{NN}}$: no informative eigenvector of the Hessian

δ_{NN} locates a **spectral phase transition**.

⁴Non-vanishing correlation with the ‘non-trivial’ directions of Θ_* .

Our results (informal overview)

This spectral phase transition threshold (BBP threshold) is well-conjectured to be the threshold for weak recovery⁵ of Θ_* via gradient descent:

▶ $\delta > \delta_{\text{NN}} \implies$ Weak recovery of Θ_*

▶ $\delta < \delta_{\text{NN}} \implies$ Failure of weak recovery under GD

⁵Non-vanishing correlation between the ‘non-trivial’ directions of Θ_* and the learned weights Θ .

Properties of δ_{NN}

① **Computable:**

δ_{NN} can be efficiently computed.

② **Predictive:**

δ_{NN} well-predicts the feature learning phase transition in numerical simulations.

③ **Sub-optimal:** $\delta_{\text{alg}} < \delta_{\text{NN}}$

The gap depends on activation, loss, initialization, ...

Properties of δ_{NN}

① **Computable:**

δ_{NN} can be efficiently computed.

② **Predictive:**

δ_{NN} well-predicts the feature learning phase transition in numerical simulations.

③ **Sub-optimal:** $\delta_{\text{alg}} < \delta_{\text{NN}}$

The gap depends on activation, loss, initialization, ...

Properties of δ_{NN}

① **Computable:**

δ_{NN} can be efficiently computed.

② **Predictive:**

δ_{NN} well-predicts the feature learning phase transition in numerical simulations.

③ **Sub-optimal:** $\delta_{\text{alg}} < \delta_{\text{NN}}$

The gap depends on activation, loss, initialization, ...

Grokking phenomenon

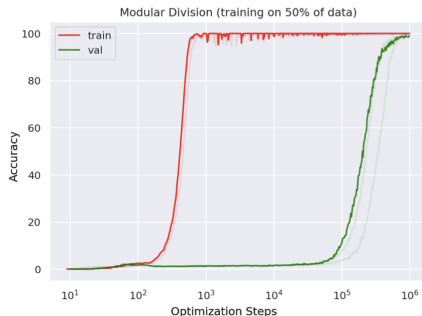
Grokking: Delayed generalization

- ▶ Training error drops quickly
- ▶ Test error stays high for long time
- ▶ **Suddenly:** test error drops

Example: Modular division task, trained on transformer architecture.

- ▶ **Red:** train accuracy
- ▶ **Green:** validation accuracy
- ▶ Gap of $\sim 10^3 \times$ in steps

[Power, Burda, Edwards+ '22], [Liu, Kitouni, Nolte+ '22], ...



Grokking phenomenon

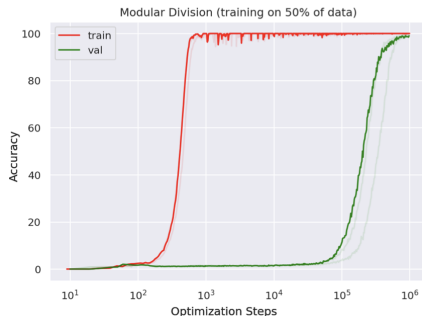
Grokking: Delayed generalization

- ▶ Training error drops quickly
- ▶ Test error stays high for long time
- ▶ **Suddenly:** test error drops

Example: Modular division task, trained on transformer architecture.

- ▶ **Red:** train accuracy
- ▶ **Green:** validation accuracy
- ▶ Gap of $\sim 10^3 \times$ in steps

[Power, Burda, Edwards+ '22], [Liu, Kitoui, Nolte+ '22], ...



Grokking phenomenon

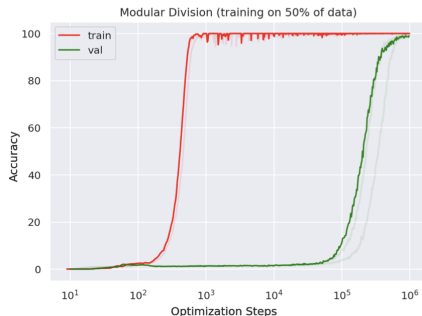
Grokking: Delayed generalization

- ▶ Training error drops quickly
- ▶ Test error stays high for long time
- ▶ **Suddenly:** test error drops

Example: Modular division task, trained on transformer architecture.

- ▶ **Red:** train accuracy
- ▶ **Green:** validation accuracy
- ▶ Gap of $\sim 10^3 \times$ in steps

[Power, Burda, Edwards+ '22], [Liu, Kitoui, Nolte+ '22], ...



Grokking phenomenon

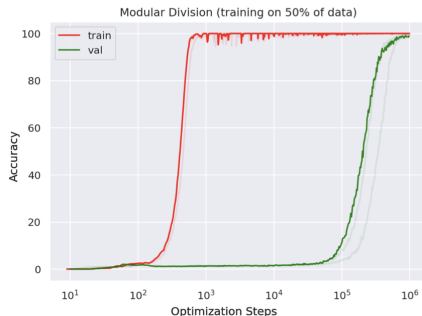
Grokking: Delayed generalization

- ▶ Training error drops quickly
- ▶ Test error stays high for long time
- ▶ **Suddenly:** test error drops

Example: Modular division task, trained on transformer architecture.

- ▶ **Red:** train accuracy
- ▶ **Green:** validation accuracy
- ▶ Gap of $\sim 10^3 \times$ in steps

[Power, Burda, Edwards+ '22], [Liu, Kitoui, Nolte+ '22], ...



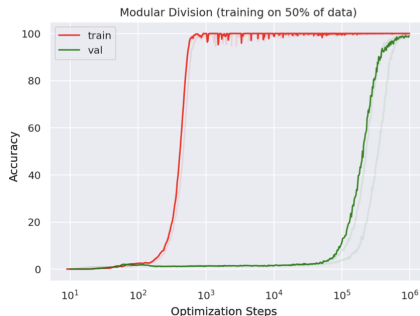
Grokking phenomenon

Grokking: Delayed generalization

- ▶ Training error drops quickly
- ▶ Test error stays high for long time
- ▶ **Suddenly:** test error drops

Example: Modular division task, trained on transformer architecture.

- ▶ **Red:** train accuracy
- ▶ **Green:** validation accuracy
- ▶ Gap of $\sim 10^3 \times$ in steps



[Power, Burda, Edwards+ '22], [Liu, Kitoui, Nolte+ '22], ...

Grokking phenomenon

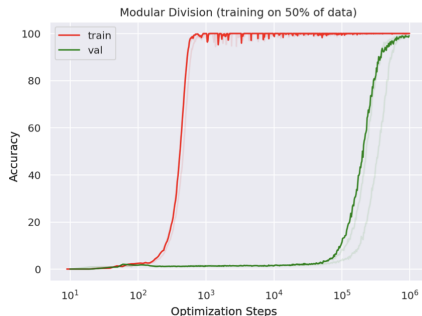
Grokking: Delayed generalization

- ▶ Training error drops quickly
- ▶ Test error stays high for long time
- ▶ **Suddenly:** test error drops

Example: Modular division task, trained on transformer architecture.

- ▶ **Red:** train accuracy
- ▶ **Green:** validation accuracy
- ▶ Gap of $\sim 10^3 \times$ in steps

[Power, Burda, Edwards+ '22], [Liu, Kitouni, Nolte+ '22], ...



Demystifying grokking

Insights:

Grokking occurs
when and only when δ is **moderately above** δ_{NN}

▶ **Moderately above threshold:**

- ▶ Overfit the training data in the first stage of training
- ▶ Invertive lesson order in the second stage of training → Grokking

▶ **Far above threshold ($\delta \gg \delta_{NN}$):**

- ▶ Language generalizes to generalization gap

Demystifying grokking

Insights:

Grokking occurs
when and only when δ is **moderately above** δ_{NN}

- ▶ **Moderately above threshold:**
 - ▶ Overfit the training data in the first stage of training
 - ▶ Informative Hessian outlier in the second stage of training \Rightarrow
Grokking
- ▶ **Far above threshold** ($\delta \gg \delta_{\text{NN}}$):
 - ▶ Landscape concentrates, no generalization gap

Demystifying grokking

Insights:

Grokking occurs
when and only when δ is **moderately above** δ_{NN}

- ▶ **Moderately above threshold:**
 - ▶ Overfit the training data in the first stage of training
 - ▶ Informative Hessian outlier in the second stage of training \Rightarrow
Grokking
- ▶ **Far above threshold** ($\delta \gg \delta_{\text{NN}}$):
 - ▶ Landscape concentrates, no generalization gap

Demystifying grokking

Insights:

Grokking occurs
when and only when δ is **moderately above** δ_{NN}

- ▶ **Moderately above threshold:**
 - ▶ Overfit the training data in the first stage of training
 - ▶ Informative Hessian outlier in the second stage of training \Rightarrow
Grokking
- ▶ **Far above threshold** ($\delta \gg \delta_{\text{NN}}$):
 - ▶ Landscape concentrates, no generalization gap

Demystifying grokking

Insights:

Grokking occurs
when and only when δ is **moderately above** δ_{NN}

- ▶ **Moderately above threshold:**
 - ▶ Overfit the training data in the first stage of training
 - ▶ Informative Hessian outlier in the second stage of training \Rightarrow
Grokking
- ▶ **Far above threshold** ($\delta \gg \delta_{\text{NN}}$):
 - ▶ Landscape concentrates, no generalization gap

Demystifying grokking

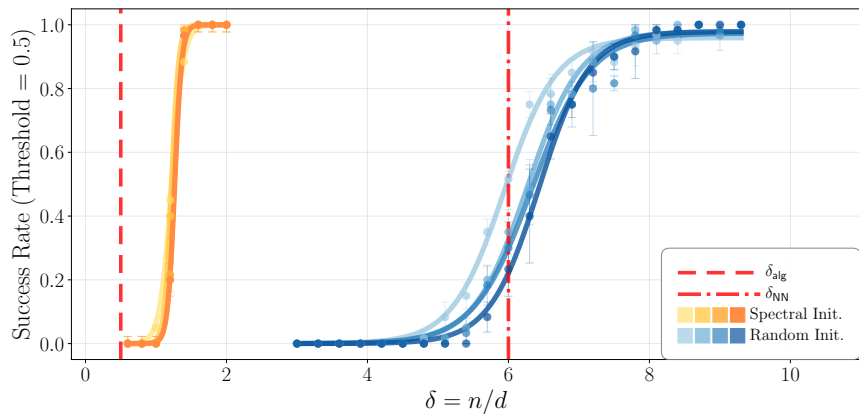
Insights:

Grokking occurs
when and only when δ is **moderately above** δ_{NN}

- ▶ **Moderately above threshold:**
 - ▶ Overfit the training data in the first stage of training
 - ▶ Informative Hessian outlier in the second stage of training \Rightarrow
Grokking
- ▶ **Far above threshold** ($\delta \gg \delta_{\text{NN}}$):
 - ▶ Landscape concentrates, no generalization gap

Numerical Illustrations

Experiments: Phase transition for feature learning

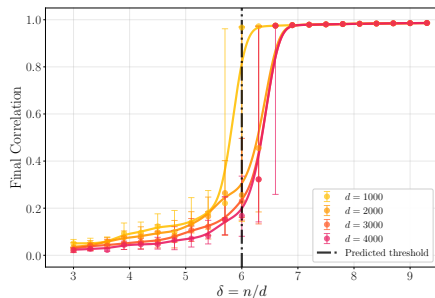


Single-neuron GeLU network learning phase retrieval ($\mathbf{h}(z) = z^2$).

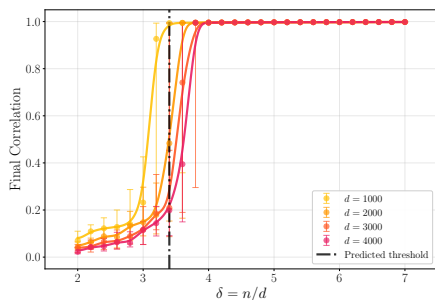
Success = correlation with $\boldsymbol{\theta}_*$ exceeds $1/2$.

Effect of activations

GeLU activation



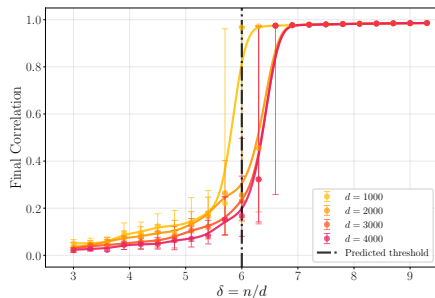
Quadratic activation



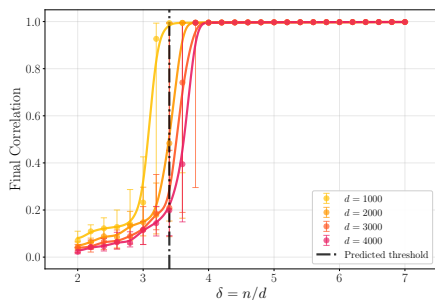
- ▶ Different activations \Rightarrow different computed δ_{NN}
- ▶ Theoretical threshold matches empirical phase transition

Effect of activations

GeLU activation

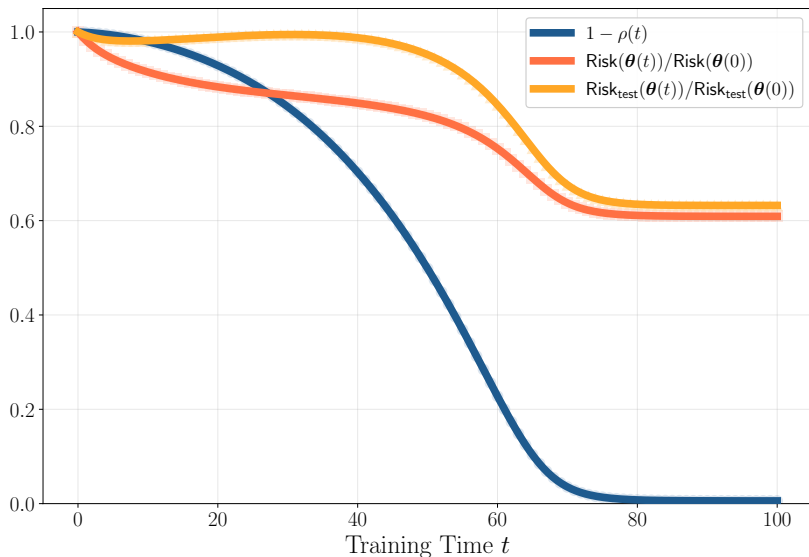


Quadratic activation



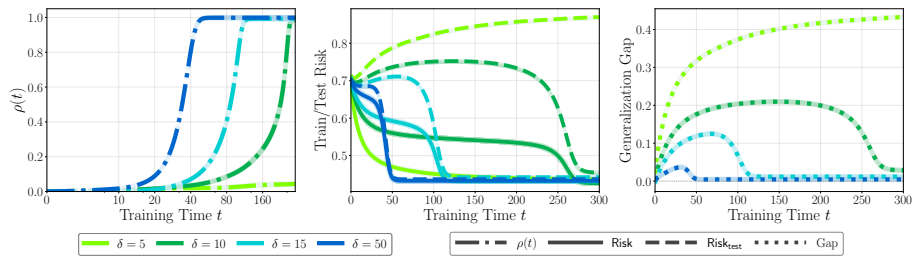
- ▶ Different activations \Rightarrow different computed δ_{NN}
- ▶ Theoretical threshold matches empirical phase transition

Experiment: Grokking phenomenon



GeLU neuron, phase retrieval, $d = 5000$, $\delta = 17.5$ (above $\delta_{\text{NN}} \approx 6$).

Experiment: Grokking



GeLU neuron, phase retrieval, $d = 5000$, $\delta = 5, 10, 15, 50$.

Formal results and Outline of the Proofs

Structural decomposition: Insight

The latent space $\text{span}(\Theta_*)$ decomposes into two parts:

Easy subspace

- ▶ Learned from **first-order** info of $h(\cdot)$
- ▶ Gradient descent can pick up **linear signal**

Hard subspace

- ▶ **No first-order** info from $h(\cdot)$
- ▶ Requires at least **2nd-order** info to learn!

Structural decomposition: Insight

The latent space $\text{span}(\Theta_*)$ decomposes into two parts:

Easy subspace

- ▶ Learned from **first-order** info of $h(\cdot)$
- ▶ Gradient descent can pick up **linear signal**

Hard subspace

- ▶ **No first-order** info from $h(\cdot)$
- ▶ Requires at least **2nd-order** info to learn!

Structural decomposition: Insight

The latent space $\text{span}(\Theta_*)$ decomposes into two parts:

Easy subspace

- ▶ Learned from **first-order** info of $h(\cdot)$
- ▶ Gradient descent can pick up **linear signal**

Hard subspace

- ▶ **No first-order** info from $h(\cdot)$
- ▶ Requires at least **2nd-order** info to learn!

Structural decomposition: Formal definition

Hard subspace: $\mathbf{U} \in \mathcal{O}(k, k')$ is a hard orthonormal set if

$$\mathbb{E}[\mathcal{T}(\mathbf{y}, \mathbf{P}_{\mathbf{U}}^{\perp} \mathbf{z}, \xi) \mathbf{P}_{\mathbf{U}} \mathbf{z}] = \mathbf{0}, \quad \forall \text{ measurable } \mathcal{T} : \mathbb{R}^{k+2} \rightarrow \mathbb{R}$$

where $\xi \sim \mathbf{N}(0, 1)$ is independent of (\mathbf{y}, \mathbf{z}) (recall $\mathbf{y} = \mathbf{h}(\mathbf{z}, \varepsilon)$).

The hard subspace $\mathfrak{U}_{\mathbf{H}} := \text{span}(\mathbf{U}_{\mathbf{H}}^*)$ is the span of the **maximal hard orthonormal set**.

$$\mathfrak{U}_{\mathbf{E}} := \mathfrak{U}_{\mathbf{H}}^{\perp}$$

Equivalent characterization:

$$\mathbf{U} \text{ is hard} \quad \Leftrightarrow \quad \mathbb{E}[\mathbf{P}_{\mathbf{U}} \mathbf{z} \mid \mathbf{y}, \mathbf{P}_{\mathbf{U}}^{\perp} \mathbf{z}] = \mathbf{0}$$

Structural decomposition: Formal definition

Hard subspace: $\mathbf{U} \in \mathcal{O}(k, k')$ is a hard orthonormal set if

$$\mathbb{E}[\mathcal{T}(\mathbf{y}, \mathbf{P}_{\mathbf{U}}^{\perp} \mathbf{z}, \xi) \mathbf{P}_{\mathbf{U}} \mathbf{z}] = \mathbf{0}, \quad \forall \text{ measurable } \mathcal{T} : \mathbb{R}^{k+2} \rightarrow \mathbb{R}$$

where $\xi \sim \mathbf{N}(0, 1)$ is independent of (\mathbf{y}, \mathbf{z}) (recall $\mathbf{y} = \mathbf{h}(\mathbf{z}, \varepsilon)$).

The hard subspace $\mathfrak{U}_{\mathbf{H}} := \text{span}(\mathbf{U}_{\mathbf{H}}^*)$ is the span of the **maximal hard orthonormal set**.

$$\mathfrak{U}_{\mathbf{E}} := \mathfrak{U}_{\mathbf{H}}^{\perp}$$

Equivalent characterization:

$$\mathbf{U} \text{ is hard} \quad \Leftrightarrow \quad \mathbb{E}[\mathbf{P}_{\mathbf{U}} \mathbf{z} \mid \mathbf{y}, \mathbf{P}_{\mathbf{U}}^{\perp} \mathbf{z}] = \mathbf{0}$$

Structural decomposition: Formal definition

Hard subspace: $\mathbf{U} \in \mathcal{O}(k, k')$ is a hard orthonormal set if

$$\mathbb{E}[\mathcal{T}(\mathbf{y}, \mathbf{P}_{\mathbf{U}}^{\perp} \mathbf{z}, \xi) \mathbf{P}_{\mathbf{U}} \mathbf{z}] = \mathbf{0}, \quad \forall \text{ measurable } \mathcal{T} : \mathbb{R}^{k+2} \rightarrow \mathbb{R}$$

where $\xi \sim \mathcal{N}(0, 1)$ is independent of (\mathbf{y}, \mathbf{z}) (recall $\mathbf{y} = \mathbf{h}(\mathbf{z}, \varepsilon)$).

The hard subspace $\mathfrak{U}_{\mathbf{H}} := \text{span}(\mathbf{U}_{\mathbf{H}}^*)$ is the span of the **maximal hard orthonormal set**.

$$\mathfrak{U}_{\mathbf{E}} := \mathfrak{U}_{\mathbf{H}}^{\perp}$$

Equivalent characterization:

$$\mathbf{U} \text{ is hard} \quad \Leftrightarrow \quad \mathbb{E}[\mathbf{P}_{\mathbf{U}} \mathbf{z} \mid \mathbf{y}, \mathbf{P}_{\mathbf{U}}^{\perp} \mathbf{z}] = \mathbf{0}$$

Structural decomposition: Formal definition

Hard subspace: $\mathbf{U} \in \mathcal{O}(k, k')$ is a hard orthonormal set if

$$\mathbb{E}[\mathcal{T}(\mathbf{y}, \mathbf{P}_{\mathbf{U}}^{\perp} \mathbf{z}, \xi) \mathbf{P}_{\mathbf{U}} \mathbf{z}] = \mathbf{0}, \quad \forall \text{ measurable } \mathcal{T} : \mathbb{R}^{k+2} \rightarrow \mathbb{R}$$

where $\xi \sim \mathcal{N}(0, 1)$ is independent of (\mathbf{y}, \mathbf{z}) (recall $\mathbf{y} = \mathbf{h}(\mathbf{z}, \varepsilon)$).

The hard subspace $\mathfrak{U}_{\mathbf{H}} := \text{span}(\mathbf{U}_{\mathbf{H}}^*)$ is the span of the **maximal hard orthonormal set**.

$$\mathfrak{U}_{\mathbf{E}} := \mathfrak{U}_{\mathbf{H}}^{\perp}$$

Equivalent characterization:

$$\mathbf{U} \text{ is hard} \quad \Leftrightarrow \quad \mathbb{E}[\mathbf{P}_{\mathbf{U}} \mathbf{z} \mid \mathbf{y}, \mathbf{P}_{\mathbf{U}}^{\perp} \mathbf{z}] = \mathbf{0}$$

Examples

Setup: $k = 2$, $y = h(z_1, z_2)$ with h symmetric: $h(z_1, z_2) = h(z_2, z_1)$

Decomposition:

- ▶ **Easy direction:** $\mathbf{u}_E = \frac{1}{\sqrt{2}}(+1, +1)^\top$
 - ▶ $\mathbb{E}[z_1 + z_2 | y] \neq 0$ in general
- ▶ **Hard direction:** $\mathbf{u}_H = \frac{1}{\sqrt{2}}(+1, -1)^\top$
 - ▶ $\mathbb{E}[z_1 - z_2 | y, z_1 + z_2] = 0$. No signal from first-order information

Other examples:

- ▶ Single-index with even $h \Rightarrow$ the only direction is hard
- ▶ Two-layer NN with even activation \Rightarrow all directions are hard

Examples

Setup: $k = 2$, $y = h(z_1, z_2)$ with h symmetric: $h(z_1, z_2) = h(z_2, z_1)$

Decomposition:

- ▶ **Easy direction:** $\mathbf{u}_E = \frac{1}{\sqrt{2}}(+1, +1)^\top$
 - ▶ $\mathbb{E}[z_1 + z_2 \mid y] \neq 0$ in general
- ▶ **Hard direction:** $\mathbf{u}_H = \frac{1}{\sqrt{2}}(+1, -1)^\top$
 - ▶ $\mathbb{E}[z_1 - z_2 \mid y, z_1 + z_2] = 0$. No signal from first-order information

Other examples:

- ▶ Single-index with even $h \Rightarrow$ the only direction is hard
- ▶ Two-layer NN with even activation \Rightarrow all directions are hard

Examples

Setup: $k = 2$, $y = h(z_1, z_2)$ with h symmetric: $h(z_1, z_2) = h(z_2, z_1)$

Decomposition:

- ▶ **Easy direction:** $\mathbf{u}_E = \frac{1}{\sqrt{2}}(+1, +1)^\top$
 - ▶ $\mathbb{E}[z_1 + z_2 \mid y] \neq 0$ in general
- ▶ **Hard direction:** $\mathbf{u}_H = \frac{1}{\sqrt{2}}(+1, -1)^\top$
 - ▶ $\mathbb{E}[z_1 - z_2 \mid y, z_1 + z_2] = 0$. No signal from first-order information

Other examples:

- ▶ Single-index with even $h \Rightarrow$ the only direction is hard
- ▶ Two-layer NN with even activation \Rightarrow all directions are hard

Examples

Setup: $k = 2$, $y = h(z_1, z_2)$ with h symmetric: $h(z_1, z_2) = h(z_2, z_1)$

Decomposition:

- ▶ **Easy direction:** $\mathbf{u}_E = \frac{1}{\sqrt{2}}(+1, +1)^\top$
 - ▶ $\mathbb{E}[z_1 + z_2 \mid y] \neq 0$ in general
- ▶ **Hard direction:** $\mathbf{u}_H = \frac{1}{\sqrt{2}}(+1, -1)^\top$
 - ▶ $\mathbb{E}[z_1 - z_2 \mid y, z_1 + z_2] = 0$. No signal from first-order information

Other examples:

- ▶ Single-index with even $h \Rightarrow$ the only direction is hard
- ▶ Two-layer NN with even activation \Rightarrow all directions are hard

Hard directions: Not learned in $O(1)$ time

Result: For any fixed t

$$\text{p-lim}_{n,d \rightarrow \infty} \Theta(t)^\top \Theta_* \mathbf{P}_{\mathbf{u}_H^*} = 0$$

More general: Any estimator computed from $O(1)$ gradient steps is uncorrelated with hard directions.

Tool: Dynamical Mean Field Theory (DMFT)⁶

⁶[Sompolinsky, Zippelius '81, '82; Celentano, Montanari, Wu '20; Celentano, Cheng, Montanari '21]

Hard directions: Not learned in $O(1)$ time

Result: For any fixed t

$$\text{p-lim}_{n,d \rightarrow \infty} \Theta(t)^\top \Theta_* P_{U_H^*} = 0$$

More general: Any estimator computed from $O(1)$ gradient steps is uncorrelated with hard directions.

Tool: Dynamical Mean Field Theory (DMFT)⁶

⁶[Sompolinsky, Zippelius '81, '82; Celentano, Montanari, Wu '20; Celentano, Cheng, Montanari '21]

Hard directions: Not learned in $O(1)$ time

Result: For any fixed t

$$\text{p-lim}_{n,d \rightarrow \infty} \Theta(t)^\top \Theta_* P_{U_H^*} = 0$$

More general: Any estimator computed from $O(1)$ gradient steps is uncorrelated with hard directions.

Tool: Dynamical Mean Field Theory (DMFT)⁶

⁶[Sompolinsky, Zippelius '81, '82; Celentano, Montanari, Wu '20; Celentano, Cheng, Montanari '21]

Easy directions: Learned in $O(1)$ time

There exists $\varepsilon_0 > 0$, $t = O(1)$ such that⁷

$$\lim_{n, d \rightarrow \infty} \mathbb{P} \left(\left\| \Theta(t)^\top \Theta_* \mathbf{u}_E^* \right\| \geq \varepsilon_0 \right) = 1$$

Intuition:

- ▶ Easy directions have **first-order signal** in the gradient

The learning bottleneck is the hard directions.

⁷[Dandi, Troiani, Arnaboldi, Pesce, Zdeborová, Krzakala '24]

Easy directions: Learned in $O(1)$ time

There exists $\varepsilon_0 > 0$, $t = O(1)$ such that⁷

$$\lim_{n, d \rightarrow \infty} \mathbb{P} \left(\left\| \Theta(t)^\top \Theta_* \mathbf{u}_E^* \right\| \geq \varepsilon_0 \right) = 1$$

Intuition:

- ▶ Easy directions have **first-order signal** in the gradient

The learning bottleneck is the hard directions.

⁷[Dandi, Troiani, Arnaboldi, Pesce, Zdeborová, Krzakala '24]

Easy directions: Learned in $O(1)$ time

There exists $\varepsilon_0 > 0$, $t = O(1)$ such that⁷

$$\lim_{n, d \rightarrow \infty} \mathbb{P} \left(\left\| \Theta(t)^\top \Theta_* \mathbf{u}_E^* \right\| \geq \varepsilon_0 \right) = 1$$

Intuition:

- ▶ Easy directions have **first-order signal** in the gradient

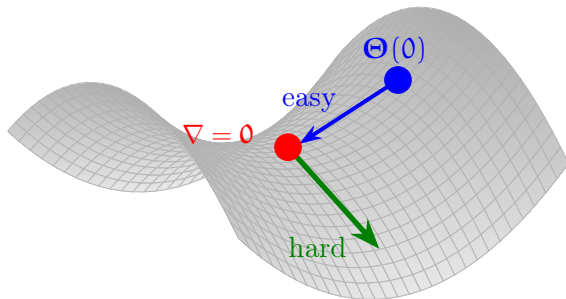
The learning bottleneck is the hard directions.

⁷[Dandi, Troiani, Arnaboldi, Pesce, Zdeborová, Krzakala '24]

Learning hard directions: The problem

Problem: No first-order signal for hard directions

Naive, simplified picture: Gradient descent approaches a **saddle point**.



Learning hard directions: The solution

Solution: Use **second-order information to escape saddle**

- ▶ Study the (rescaled) Hessian

$$\mathbf{H}(t) := m \nabla^2 \text{Risk}(\Theta(t))$$

along GD trajectory

- ▶ Look for **negative eigenvalues** of the Hessian with eigenvectors aligned with hard subspace
- ▶ Such eigenvectors directions enable **escape from saddle** \Rightarrow learning those hard directions

Question: When does the Hessian have **informative descent directions** after $O(1)$ iterations?

Learning hard directions: The solution

Solution: Use **second-order information to escape saddle**

- ▶ Study the (rescaled) Hessian

$$\mathbf{H}(t) := m \nabla^2 \text{Risk}(\Theta(t))$$

along GD trajectory

- ▶ Look for **negative eigenvalues** of the Hessian with eigenvectors aligned with hard subspace
- ▶ Such eigenvectors directions enable **escape from saddle** \Rightarrow learning those hard directions

Question: When does the Hessian have **informative descent directions** after $O(1)$ iterations?

Learning hard directions: The solution

Solution: Use **second-order information to escape saddle**

- ▶ Study the (rescaled) Hessian

$$\mathbf{H}(t) := m \nabla^2 \text{Risk}(\Theta(t))$$

along GD trajectory

- ▶ Look for **negative eigenvalues** of the Hessian with eigenvectors aligned with hard subspace
- ▶ Such eigenvectors directions enable **escape from saddle** \Rightarrow learning those hard directions

Question: When does the Hessian have **informative descent directions** after $O(1)$ iterations?

Learning hard directions: The solution

Solution: Use **second-order information to escape saddle**

- ▶ Study the (rescaled) Hessian

$$\mathbf{H}(t) := m \nabla^2 \text{Risk}(\Theta(t))$$

along GD trajectory

- ▶ Look for **negative eigenvalues** of the Hessian with eigenvectors aligned with hard subspace
- ▶ Such eigenvectors directions enable **escape from saddle** \Rightarrow learning those hard directions

Question: When does the Hessian have **informative descent directions** after $O(1)$ iterations?

Learning hard directions: The solution

Solution: Use **second-order information to escape saddle**

- ▶ Study the (rescaled) Hessian

$$\mathbf{H}(t) := m \nabla^2 \text{Risk}(\Theta(t))$$

along GD trajectory

- ▶ Look for **negative eigenvalues** of the Hessian with eigenvectors aligned with hard subspace
- ▶ Such eigenvectors directions enable **escape from saddle** \Rightarrow learning those hard directions

Question: When does the Hessian have **informative descent directions** after $O(1)$ iterations?

Hessian structure

Case $m = 1$ (single neuron): $\mathbf{H}(t) \in \mathbb{R}^{d \times d}$

$$\mathbf{H}(t) = \frac{1}{n} \sum_{i=1}^n g(y_i, \boldsymbol{\theta}(t)^\top \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top$$

$$g(y, z) = a^2 \ell''(y, a\sigma(z + b)) \sigma'(z + b)^2 + a \ell'(y, a\sigma(z + b)) \sigma''(z + b)$$

Hessian structure

Case $m \gg 1$ (wide network): $\mathbf{H}(\mathbf{t}) \in \mathbb{R}^{md \times md}$

- ▶ Block structure with m^2 blocks of size $d \times d$
- ▶ For convex loss: well-approximated by block-diagonal $\mathbf{H}_{\text{diag}}(\mathbf{t})$
- ▶ Each block $\mathbf{a}_j \mathbf{H}_j(\mathbf{t})$ where

$$\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \boldsymbol{\Theta}(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top$$
$$g(\mathbf{y}, \boldsymbol{\Theta}^\top \mathbf{x}; j) = \ell'(\mathbf{y}, f_{\boldsymbol{\Theta}}(\mathbf{x})) \cdot \sigma''(\boldsymbol{\theta}_j^\top \mathbf{x} + \mathbf{b}_j)$$

Both cases reduce to **spiked random matrix** with training-dependent weights.

Hessian structure

Case $m \gg 1$ (wide network): $\mathbf{H}(\mathbf{t}) \in \mathbb{R}^{md \times md}$

- ▶ Block structure with m^2 blocks of size $d \times d$
- ▶ For convex loss: well-approximated by block-diagonal $\mathbf{H}_{\text{diag}}(\mathbf{t})$
- ▶ Each block $\mathbf{a}_j \mathbf{H}_j(\mathbf{t})$ where

$$\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \Theta(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top$$
$$g(\mathbf{y}, \Theta^\top \mathbf{x}; j) = \ell'(\mathbf{y}, f_{\Theta}(\mathbf{x})) \cdot \sigma''(\theta_j^\top \mathbf{x} + \mathbf{b}_j)$$

Both cases reduce to **spiked random matrix** with training-dependent weights.

Main theorem: Hessian phase transition

(Stated for $m \gg 1$; similar theorem holds for $m = 1$)

Theorem: Under $n, d \rightarrow \infty$ with $n/d \rightarrow \delta$, under generic regularity assumptions on \mathbf{h}, ℓ, σ , there exist explicit thresholds $\delta_j^* := \mathcal{K}_j(\mathbf{h}, \ell, \sigma, \eta)$ for each block j :

If $\delta > \delta_j^*$: For t a large enough constant, there are $\lambda_{n,d}, \xi_{n,d}$ such that $\mathbf{H}_j(t)\xi_{n,d} = \lambda_{n,d}\xi_{n,d}$, and

$$\lambda_{n,d} \xrightarrow{P} \lambda_j^* < \min(0, \inf(\text{supp}(\mu_\infty^j(t))))$$
$$\left\| \Theta_{*H}^\top \xi_{n,d} \right\| / \|\xi_{n,d}\| \xrightarrow{P} c_j^* > 0$$

where $\mu_\infty^j(t)$ is the limiting spectral distribution of $\mathbf{H}_j(t)$.

If $\delta < \delta_j^*$: For all large enough t a constant, no eigenvector of $\mathbf{H}_j(t)$ has non-vanishing correlation with the hard subspace

$\Rightarrow \delta_{NN} := \min_j \delta_j^*$ is the threshold for feature learning!

Main theorem: Hessian phase transition

(Stated for $m \gg 1$; similar theorem holds for $m = 1$)

Theorem: Under $n, d \rightarrow \infty$ with $n/d \rightarrow \delta$, under generic regularity assumptions on \mathbf{h}, ℓ, σ , there exist explicit thresholds $\delta_j^* := \mathcal{K}_j(\mathbf{h}, \ell, \sigma, \eta)$ for each block j :

If $\delta > \delta_j^*$: For t a large enough constant, there are $\lambda_{n,d}, \xi_{n,d}$ such that $\mathbf{H}_j(t)\xi_{n,d} = \lambda_{n,d}\xi_{n,d}$, and

$$\lambda_{n,d} \xrightarrow{P} \lambda_j^* < \min(0, \inf(\text{supp}(\mu_\infty^j(t))))$$
$$\left\| \Theta_{*H}^\top \xi_{n,d} \right\| / \|\xi_{n,d}\| \xrightarrow{P} c_j^* > 0$$

where $\mu_\infty^j(t)$ is the limiting spectral distribution of $\mathbf{H}_j(t)$.

If $\delta < \delta_j^*$: For all large enough t a constant, no eigenvector of $\mathbf{H}_j(t)$ has non-vanishing correlation with the hard subspace

$\Rightarrow \delta_{NN} := \min_j \delta_j^*$ is the threshold for feature learning!

Main theorem: Hessian phase transition

(Stated for $m \gg 1$; similar theorem holds for $m = 1$)

Theorem: Under $n, d \rightarrow \infty$ with $n/d \rightarrow \delta$, under generic regularity assumptions on \mathbf{h}, ℓ, σ , there exist explicit thresholds $\delta_j^* := \mathcal{K}_j(\mathbf{h}, \ell, \sigma, \eta)$ for each block j :

If $\delta > \delta_j^*$: For t a large enough constant, there are $\lambda_{n,d}, \xi_{n,d}$ such that $\mathbf{H}_j(t)\xi_{n,d} = \lambda_{n,d}\xi_{n,d}$, and

$$\lambda_{n,d} \xrightarrow{P} \lambda_j^* < \min(0, \inf(\text{supp}(\mu_\infty^j(t))))$$
$$\left\| \Theta_{*H}^\top \xi_{n,d} \right\| / \|\xi_{n,d}\| \xrightarrow{P} c_j^* > 0$$

where $\mu_\infty^j(t)$ is the limiting spectral distribution of $\mathbf{H}_j(t)$.

If $\delta < \delta_j^*$: For all large enough t a constant, no eigenvector of $\mathbf{H}_j(t)$ has non-vanishing correlation with the hard subspace

$\Rightarrow \delta_{NN} := \min_j \delta_j^*$ is the threshold for feature learning!

Main theorem: Hessian phase transition

(Stated for $m \gg 1$; similar theorem holds for $m = 1$)

Theorem: Under $n, d \rightarrow \infty$ with $n/d \rightarrow \delta$, under generic regularity assumptions on \mathbf{h}, ℓ, σ , there exist explicit thresholds $\delta_j^* := \mathcal{K}_j(\mathbf{h}, \ell, \sigma, \eta)$ for each block j :

If $\delta > \delta_j^*$: For t a large enough constant, there are $\lambda_{n,d}, \xi_{n,d}$ such that $\mathbf{H}_j(t)\xi_{n,d} = \lambda_{n,d}\xi_{n,d}$, and

$$\lambda_{n,d} \xrightarrow{P} \lambda_j^* < \min(0, \inf(\text{supp}(\mu_\infty^j(t))))$$
$$\left\| \Theta_{*H}^\top \xi_{n,d} \right\| / \|\xi_{n,d}\| \xrightarrow{P} c_j^* > 0$$

where $\mu_\infty^j(t)$ is the limiting spectral distribution of $\mathbf{H}_j(t)$.

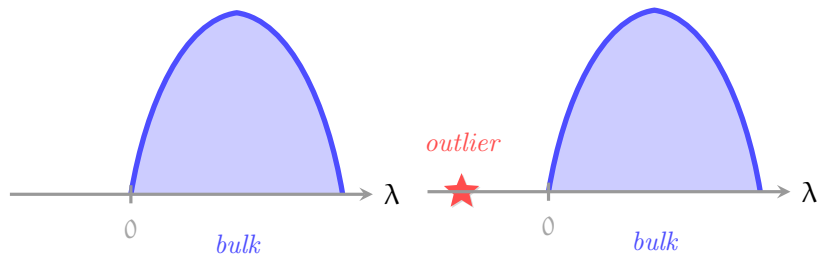
If $\delta < \delta_j^*$: For all large enough t a constant, no eigenvector of $\mathbf{H}_j(t)$ has non-vanishing correlation with the hard subspace

$\Rightarrow \delta_{NN} := \min_j \delta_j^*$ is the threshold for feature learning!

Spectrum of the Hessian

$$\delta < \delta_{NN}$$

$$\delta > \delta_{NN}$$



No informative direction

Descent direction \rightarrow hard subspace!

Formulas for computing δ_{NN}

Hessian block: $\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \boldsymbol{\Theta}(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top$

Limiting spectral distribution: $\mu_{n,d}^j(\mathbf{t}) \Rightarrow \mu_\infty^j(\mathbf{t})$ with $\alpha_t^j(z)$ the Stieltjes transform of $\mu_\infty^j(\mathbf{t})$.

Stieltjes transform: For $z \in \mathbb{C}^+$, $\alpha_t^j(z)$ is the unique solution of

$$z + \frac{1}{\alpha_t^j(z)} = \delta \cdot \mathbb{E} \left[\frac{G_t^j}{\delta + G_t^j \alpha_t^j(z)} \right]$$

where $G_t^j := g(\mathbf{V}(\mathbf{t}), \mathbf{h}(\mathbf{V}_*, \varepsilon); j)$.

$(\mathbf{V}(\mathbf{t}), \mathbf{V}_*)$ is some DMFT random vectors that we define later. The resulting LSD is a **Generalized Marchenko-Pastur law** despite complex dependence.

Formulas for computing δ_{NN}

Hessian block: $\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \boldsymbol{\Theta}(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top$

Limiting spectral distribution: $\mu_{n,d}^j(\mathbf{t}) \Rightarrow \mu_\infty^j(\mathbf{t})$ with $\alpha_t^j(z)$ the Stieltjes transform of $\mu_\infty^j(\mathbf{t})$.

Stieltjes transform: For $z \in \mathbb{C}^+$, $\alpha_t^j(z)$ is the unique solution of

$$z + \frac{1}{\alpha_t^j(z)} = \delta \cdot \mathbb{E} \left[\frac{G_t^j}{\delta + G_t^j \alpha_t^j(z)} \right]$$

where $G_t^j := g(\mathbf{V}(\mathbf{t}), \mathbf{h}(\mathbf{V}_*, \varepsilon); j)$.

$(\mathbf{V}(\mathbf{t}), \mathbf{V}_*)$ is some DMFT random vectors that we define later. The resulting LSD is a **Generalized Marchenko-Pastur law** despite complex dependence.

Formulas for computing δ_{NN}

Hessian block: $\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \Theta(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top$

Limiting spectral distribution: $\mu_{n,d}^j(\mathbf{t}) \Rightarrow \mu_\infty^j(\mathbf{t})$ with $\alpha_t^j(z)$ the Stieltjes transform of $\mu_\infty^j(\mathbf{t})$.

Stieltjes transform: For $z \in \mathbb{C}^+$, $\alpha_t^j(z)$ is the unique solution of

$$z + \frac{1}{\alpha_t^j(z)} = \delta \cdot \mathbb{E} \left[\frac{G_t^j}{\delta + G_t^j \alpha_t^j(z)} \right]$$

where $G_t^j := g(\mathbf{V}(\mathbf{t}), \mathbf{h}(\mathbf{V}_*, \varepsilon); j)$.

$(\mathbf{V}(\mathbf{t}), \mathbf{V}_*)$ is some DMFT random vectors that we define later. The resulting LSD is a **Generalized Marchenko-Pastur law** despite complex dependence.

Formulas for computing δ_{NN}

Hessian block: $\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \Theta(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top$

Limiting spectral distribution: $\mu_{n,d}^j(\mathbf{t}) \Rightarrow \mu_\infty^j(\mathbf{t})$ with $\alpha_t^j(z)$ the Stieltjes transform of $\mu_\infty^j(\mathbf{t})$.

Stieltjes transform: For $z \in \mathbb{C}^+$, $\alpha_t^j(z)$ is the unique solution of

$$z + \frac{1}{\alpha_t^j(z)} = \delta \cdot \mathbb{E} \left[\frac{G_t^j}{\delta + G_t^j \alpha_t^j(z)} \right]$$

where $G_t^j := g(\mathbf{V}(\mathbf{t}), \mathbf{h}(\mathbf{V}_*, \varepsilon); j)$.

$(\mathbf{V}(\mathbf{t}), \mathbf{V}_*)$ is some DMFT random vectors that we define later. The resulting LSD is a **Generalized Marchenko-Pastur law** despite complex dependence.

Formulas for computing δ_{NN}

Left edge: $A_t^j := \frac{\delta}{(-\inf(\text{supp}(\text{Law}(G_t^j)))) \vee 0}$

$$c_j(t) := \sup_{\alpha \in (0, A_t^j)} \left\{ -\frac{1}{\alpha} + \delta \cdot \mathbb{E} \left[\frac{G_t^j}{\delta + G_t^j \alpha} \right] \right\}$$

Outlier equation: For $z < \min(0, c_j(t))$,

$$\det \left(-zI_r + \mathbb{E} \left[\frac{\delta G_t^j}{\delta + G_t^j \alpha_t^j(z)} \mathbf{u}_H^{*\top} \mathbf{v}_* (\mathbf{u}_H^{*\top} \mathbf{v}_*)^\top \right] \right) = 0$$

$$\delta_j^*(t) := \inf \{ \delta > 0 : \text{outlier eq. has solution } z < \min(0, c_j(t)) \}$$

$$\delta_j^*(\infty) := \lim_{t \rightarrow \infty} \delta_j^*(t)$$

Overall: $\delta_{\text{NN}} := \min_{j \in [m]} \delta_j^*(\infty)$

Formulas for computing δ_{NN}

Left edge:
$$A_t^j := \frac{\delta}{(-\inf(\text{supp}(\text{Law}(G_t^j)))) \vee 0}$$

$$c_j(t) := \sup_{\alpha \in (0, A_t^j)} \left\{ -\frac{1}{\alpha} + \delta \cdot \mathbb{E} \left[\frac{G_t^j}{\delta + G_t^j \alpha} \right] \right\}$$

Outlier equation: For $z < \min(0, c_j(t))$,

$$\det \left(-zI_r + \mathbb{E} \left[\frac{\delta G_t^j}{\delta + G_t^j \alpha_t^j(z)} \mathbf{u}_H^{*\top} \mathbf{v}_* (\mathbf{u}_H^{*\top} \mathbf{v}_*)^\top \right] \right) = 0$$

$$\delta_j^*(t) := \inf \{ \delta > 0 : \text{outlier eq. has solution } z < \min(0, c_j(t)) \}$$

$$\delta_j^*(\infty) := \lim_{t \rightarrow \infty} \delta_j^*(t)$$

Overall:
$$\delta_{\text{NN}} := \min_{j \in [m]} \delta_j^*(\infty)$$

Formulas for computing δ_{NN}

Left edge:
$$A_t^j := \frac{\delta}{(-\inf(\text{supp}(\text{Law}(G_t^j)))) \vee 0}$$

$$c_j(t) := \sup_{\alpha \in (0, A_t^j)} \left\{ -\frac{1}{\alpha} + \delta \cdot \mathbb{E} \left[\frac{G_t^j}{\delta + G_t^j \alpha} \right] \right\}$$

Outlier equation: For $z < \min(0, c_j(t))$,

$$\det \left(-zI_r + \mathbb{E} \left[\frac{\delta G_t^j}{\delta + G_t^j \alpha_t^j(z)} \mathbf{u}_H^{*\top} \mathbf{v}_* (\mathbf{u}_H^{*\top} \mathbf{v}_*)^\top \right] \right) = 0$$

$$\delta_j^*(t) := \inf \{ \delta > 0 : \text{outlier eq. has solution } z < \min(0, c_j(t)) \}$$

$$\delta_j^*(\infty) := \lim_{t \rightarrow \infty} \delta_j^*(t)$$

Overall:
$$\delta_{\text{NN}} := \min_{j \in [m]} \delta_j^*(\infty)$$

Formulas for computing δ_{NN}

Left edge: $A_t^j := \frac{\delta}{(-\inf(\text{supp}(\text{Law}(G_t^j)))) \vee 0}$

$$c_j(t) := \sup_{\alpha \in (0, A_t^j)} \left\{ -\frac{1}{\alpha} + \delta \cdot \mathbb{E} \left[\frac{G_t^j}{\delta + G_t^j \alpha} \right] \right\}$$

Outlier equation: For $z < \min(0, c_j(t))$,

$$\det \left(-zI_r + \mathbb{E} \left[\frac{\delta G_t^j}{\delta + G_t^j \alpha_t^j(z)} \mathbf{u}_H^{*\top} \mathbf{v}_* (\mathbf{u}_H^{*\top} \mathbf{v}_*)^\top \right] \right) = 0$$

$$\delta_j^*(t) := \inf \{ \delta > 0 : \text{outlier eq. has solution } z < \min(0, c_j(t)) \}$$

$$\delta_j^*(\infty) := \lim_{t \rightarrow \infty} \delta_j^*(t)$$

Overall: $\delta_{\text{NN}} := \min_{j \in [m]} \delta_j^*(\infty)$

DMFT: Evolution equations

Dynamical Mean Field Theory (DMFT)⁸: Characterize gradient descent dynamics as $n, d \rightarrow \infty$ via **low-dimensional random vectors**.

$$(V(t), V_*) \in \mathbb{R}^m \times \mathbb{R}^k$$

$$V(t) = W(t) - \frac{1}{\delta} \sum_{s=0}^{t-1} R_\theta(t, s) F(V(s), V_*, \varepsilon), \quad W \sim \text{GP}(0, C_\theta)$$

$$\Theta(t) \in \mathbb{R}^m$$

$$\Theta(t+1) = \Theta(t) - \eta \sum_{s=0}^t R_\ell(t, s) \Theta(s) - \eta R_\ell(t, *) \Theta_* + \eta Q(t), \quad Q \sim \text{GP}(0, C_\ell/\delta)$$

$$\text{Gradient function: } F_j(V, V_*, \varepsilon) = (\eta/m) a_j \ell'(h(V_*, \varepsilon), f(V)) \sigma'(V_j + b_j)$$

⁸[Sompolinsky, Zippelius '81, '82; Celentano, Montanari, Wu '20; Celentano, Cheng, Montanari '21]

DMFT: Evolution equations

Dynamical Mean Field Theory (DMFT)⁸: Characterize gradient descent dynamics as $n, d \rightarrow \infty$ via **low-dimensional random vectors**.

$$(V(t), V_*) \in \mathbb{R}^m \times \mathbb{R}^k$$

$$V(t) = W(t) - \frac{1}{\delta} \sum_{s=0}^{t-1} R_\theta(t, s) F(V(s), V_*, \varepsilon), \quad W \sim \text{GP}(0, C_\theta)$$

$$\Theta(t) \in \mathbb{R}^m$$

$$\Theta(t+1) = \Theta(t) - \eta \sum_{s=0}^t R_\ell(t, s) \Theta(s) - \eta R_\ell(t, *) \Theta_* + \eta Q(t), \quad Q \sim \text{GP}(0, C_\ell/\delta)$$

$$\text{Gradient function: } F_j(V, V_*, \varepsilon) = (\eta/m) a_j \ell'(h(V_*, \varepsilon), f(V)) \sigma'(V_j + b_j)$$

⁸[Sompolinsky, Zippelius '81, '82; Celentano, Montanari, Wu '20; Celentano, Cheng, Montanari '21]

DMFT: Evolution equations

Dynamical Mean Field Theory (DMFT)⁸: Characterize gradient descent dynamics as $n, d \rightarrow \infty$ via **low-dimensional random vectors**.

$$(V(t), V_*) \in \mathbb{R}^m \times \mathbb{R}^k$$

$$V(t) = W(t) - \frac{1}{\delta} \sum_{s=0}^{t-1} R_{\theta}(t, s) F(V(s), V_*, \varepsilon), \quad W \sim \text{GP}(0, C_{\theta})$$

$$\Theta(t) \in \mathbb{R}^m$$

$$\Theta(t+1) = \Theta(t) - \eta \sum_{s=0}^t R_{\ell}(t, s) \Theta(s) - \eta R_{\ell}(t, *) \Theta_* + \eta Q(t), \quad Q \sim \text{GP}(0, C_{\ell}/\delta)$$

Gradient function: $F_j(V, V_*, \varepsilon) = (\eta/m) a_j \ell'(h(V_*, \varepsilon), f(V)) \sigma'(V_j + b_j)$

⁸[Sompolinsky, Zippelius '81, '82; Celentano, Montanari, Wu '20; Celentano, Cheng, Montanari '21]

DMFT: Evolution equations

Dynamical Mean Field Theory (DMFT)⁸: Characterize gradient descent dynamics as $n, d \rightarrow \infty$ via **low-dimensional random vectors**.

$$(V(t), V_*) \in \mathbb{R}^m \times \mathbb{R}^k$$

$$V(t) = W(t) - \frac{1}{\delta} \sum_{s=0}^{t-1} R_\theta(t, s) F(V(s), V_*, \varepsilon), \quad W \sim \text{GP}(0, C_\theta)$$

$$\Theta(t) \in \mathbb{R}^m$$

$$\Theta(t+1) = \Theta(t) - \eta \sum_{s=0}^t R_\ell(t, s) \Theta(s) - \eta R_\ell(t, *) \Theta_* + \eta Q(t), \quad Q \sim \text{GP}(0, C_\ell/\delta)$$

Gradient function: $F_j(V, V_*, \varepsilon) = (\eta/m) a_j \ell'(h(V_*, \varepsilon), f(V)) \sigma'(V_j + b_j)$

⁸[Sompolinsky, Zippelius '81, '82; Celentano, Montanari, Wu '20; Celentano, Cheng, Montanari '21]

DMFT: Evolution equations

Covariance kernels:

$$C_{\Theta}(t, s) = \mathbb{E}[\Theta(t)\Theta(s)^{\top}], \quad C_{\Theta}(t, *) = \mathbb{E}[\Theta(t)\Theta_*^{\top}]$$

$$C_{\ell}(t, s) = \frac{1}{\eta^2} \mathbb{E}[F(V(t), V_*, \varepsilon)F(V(s), V_*, \varepsilon)^{\top}]$$

Response kernels:

$$R_{\Theta}(t, s) = \frac{1}{\eta} \mathbb{E} \left[\frac{\partial \Theta(t)}{\partial Q(s)} \right]$$

$$R_{\ell}(t, s) = \frac{1}{\eta} \mathbb{E} \left[\frac{\partial F(V(t), V_*, \varepsilon)}{\partial W(s)} \right], \quad R_{\ell}(t, *) = \frac{1}{\eta} \mathbb{E} \left[\frac{\partial F(V(t), V_*, \varepsilon)}{\partial V_*} \right]$$

Initialization:

$$\Theta(0) \sim N(0, I_m) \text{ independent of } \Theta_* \sim N(0, I_k)$$

DMFT: Evolution equations

Covariance kernels:

$$C_{\Theta}(t, s) = \mathbb{E}[\Theta(t)\Theta(s)^{\top}], \quad C_{\Theta}(t, *) = \mathbb{E}[\Theta(t)\Theta_*^{\top}]$$

$$C_{\ell}(t, s) = \frac{1}{\eta^2} \mathbb{E}[F(V(t), V_*, \varepsilon)F(V(s), V_*, \varepsilon)^{\top}]$$

Response kernels:

$$R_{\Theta}(t, s) = \frac{1}{\eta} \mathbb{E} \left[\frac{\partial \Theta(t)}{\partial Q(s)} \right]$$

$$R_{\ell}(t, s) = \frac{1}{\eta} \mathbb{E} \left[\frac{\partial F(V(t), V_*, \varepsilon)}{\partial W(s)} \right], \quad R_{\ell}(t, *) = \frac{1}{\eta} \mathbb{E} \left[\frac{\partial F(V(t), V_*, \varepsilon)}{\partial V_*} \right]$$

Initialization:

$$\Theta(0) \sim N(0, I_m) \text{ independent of } \Theta_* \sim N(0, I_k)$$

DMFT: Evolution equations

Covariance kernels:

$$C_{\Theta}(t, s) = \mathbb{E}[\Theta(t)\Theta(s)^{\top}], \quad C_{\Theta}(t, *) = \mathbb{E}[\Theta(t)\Theta_*^{\top}]$$

$$C_{\ell}(t, s) = \frac{1}{\eta^2} \mathbb{E}[F(\mathbf{V}(t), \mathbf{V}_*, \varepsilon)F(\mathbf{V}(s), \mathbf{V}_*, \varepsilon)^{\top}]$$

Response kernels:

$$R_{\Theta}(t, s) = \frac{1}{\eta} \mathbb{E} \left[\frac{\partial \Theta(t)}{\partial Q(s)} \right]$$

$$R_{\ell}(t, s) = \frac{1}{\eta} \mathbb{E} \left[\frac{\partial F(\mathbf{V}(t), \mathbf{V}_*, \varepsilon)}{\partial W(s)} \right], \quad R_{\ell}(t, *) = \frac{1}{\eta} \mathbb{E} \left[\frac{\partial F(\mathbf{V}(t), \mathbf{V}_*, \varepsilon)}{\partial \mathbf{V}_*} \right]$$

Initialization:

$$\Theta(0) \sim \mathbf{N}(0, \mathbf{I}_m) \text{ independent of } \Theta_* \sim \mathbf{N}(0, \mathbf{I}_k)$$

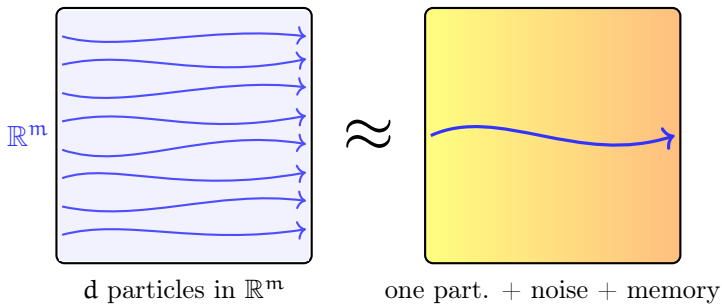
Justification of DMFT

For test function ϕ ,

$$\frac{1}{n} \sum_{i=1}^n \phi \left(\Theta(t)^\top \mathbf{x}_i, \Theta_*^\top \mathbf{x}_i, \varepsilon_i \right) \xrightarrow{P} \mathbb{E}[\phi(\mathbf{V}(t), \mathbf{V}_*, \varepsilon)],$$

$$\frac{1}{d} \sum_{i=1}^d \phi \left(\sqrt{d} \Theta(t)_i, \sqrt{d} \Theta_{*i} \right) \xrightarrow{P} \mathbb{E}[\phi(\Theta(t), \Theta_*)].$$

DMFT: Main idea



'Gaussian' approximation + Symmetries

\Rightarrow Integral equations for memory + response kernels

Proof techniques

Our proofs utilize three main ingredients:

- ① **Dynamical Mean Field Theory (DMFT)**
 - ▶ Tracks the evolution of GD in the first stage of training
- ② **Gaussian conditioning**
- ③ **Random Matrix Theory (RMT)**
 - ▶ The later two techniques are used for characterizing the Hessian spectrum (second stage).

Proof sketch: Setup

(WLOG consider $m \gg 1$; the case $m = 1$ is similar)

Object of interest: Hessian diagonal block

$$\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \Theta(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{X}^\top \mathbf{G}_j(\mathbf{t}) \mathbf{X}$$

Intuition: \mathbf{y}_i only depends on **few directions** $\Theta_*^\top \mathbf{x}_i$
 $\Rightarrow \mathbf{G}_j(\mathbf{t})$ carries low-rank informative signal \Rightarrow expect outliers

Why is this hard? Not a standard spiked random matrix

- ▶ $\Theta(\mathbf{t})$ depends on (\mathbf{X}, \mathbf{y}) through GD updates
- ▶ Complex nonlinear dependence \Rightarrow standard RMT does not work

Key insight: Although nonlinear, each GD update is **linear in \mathbf{X}**

Proof sketch: Setup

(WLOG consider $m \gg 1$; the case $m = 1$ is similar)

Object of interest: Hessian diagonal block

$$\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \Theta(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{X}^\top \mathbf{G}_j(\mathbf{t}) \mathbf{X}$$

Intuition: \mathbf{y}_i only depends on **few directions** $\Theta_*^\top \mathbf{x}_i$
 $\Rightarrow \mathbf{G}_j(\mathbf{t})$ carries low-rank informative signal \Rightarrow expect outliers

Why is this hard? Not a standard spiked random matrix

- ▶ $\Theta(\mathbf{t})$ depends on (\mathbf{X}, \mathbf{y}) through GD updates
- ▶ Complex nonlinear dependence \Rightarrow standard RMT does not work

Key insight: Although nonlinear, each GD update is **linear in \mathbf{X}**

Proof sketch: Setup

(WLOG consider $m \gg 1$; the case $m = 1$ is similar)

Object of interest: Hessian diagonal block

$$\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \Theta(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{X}^\top \mathbf{G}_j(\mathbf{t}) \mathbf{X}$$

Intuition: \mathbf{y}_i only depends on **few directions** $\Theta_*^\top \mathbf{x}_i$
 $\Rightarrow \mathbf{G}_j(\mathbf{t})$ carries low-rank informative signal \Rightarrow expect outliers

Why is this hard? **Not** a standard spiked random matrix

- ▶ $\Theta(\mathbf{t})$ depends on (\mathbf{X}, \mathbf{y}) through GD updates
- ▶ Complex nonlinear dependence \Rightarrow standard RMT does not work

Key insight: Although nonlinear, each GD update is **linear in \mathbf{X}**

Proof sketch: Setup

(WLOG consider $m \gg 1$; the case $m = 1$ is similar)

Object of interest: Hessian diagonal block

$$\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \Theta(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{X}^\top \mathbf{G}_j(\mathbf{t}) \mathbf{X}$$

Intuition: \mathbf{y}_i only depends on **few directions** $\Theta_*^\top \mathbf{x}_i$
 $\Rightarrow \mathbf{G}_j(\mathbf{t})$ carries low-rank informative signal \Rightarrow expect outliers

Why is this hard? **Not** a standard spiked random matrix

- ▶ $\Theta(\mathbf{t})$ depends on (\mathbf{X}, \mathbf{y}) through GD updates
- ▶ Complex nonlinear dependence \Rightarrow standard RMT does not work

Key insight: Although nonlinear, each GD update is **linear in \mathbf{X}**

Proof sketch: Setup

(WLOG consider $m \gg 1$; the case $m = 1$ is similar)

Object of interest: Hessian diagonal block

$$\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \Theta(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{X}^\top \mathbf{G}_j(\mathbf{t}) \mathbf{X}$$

Intuition: \mathbf{y}_i only depends on **few directions** $\Theta_*^\top \mathbf{x}_i$

$\Rightarrow \mathbf{G}_j(\mathbf{t})$ carries low-rank informative signal \Rightarrow expect outliers

Why is this hard? **Not** a standard spiked random matrix

- ▶ $\Theta(\mathbf{t})$ depends on (\mathbf{X}, \mathbf{y}) through GD updates
- ▶ Complex nonlinear dependence \Rightarrow standard RMT does not work

Key insight: Although nonlinear, each GD update is **linear in \mathbf{X}**

Proof sketch: Setup

(WLOG consider $m \gg 1$; the case $m = 1$ is similar)

Object of interest: Hessian diagonal block

$$\mathbf{H}_j(\mathbf{t}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_i, \Theta(\mathbf{t})^\top \mathbf{x}_i; j) \cdot \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{X}^\top \mathbf{G}_j(\mathbf{t}) \mathbf{X}$$

Intuition: \mathbf{y}_i only depends on **few directions** $\Theta_*^\top \mathbf{x}_i$
 $\Rightarrow \mathbf{G}_j(\mathbf{t})$ carries low-rank informative signal \Rightarrow expect outliers

Why is this hard? **Not** a standard spiked random matrix

- ▶ $\Theta(\mathbf{t})$ depends on (\mathbf{X}, \mathbf{y}) through GD updates
- ▶ Complex nonlinear dependence \Rightarrow standard RMT does not work

Key insight: Although nonlinear, each GD update is **linear in \mathbf{X}**

Proof sketch: Step 1 – Gaussian conditioning

Gaussian conditioning decomposition:

$$\mathbf{X} \stackrel{d}{=} \mathbf{P}_{\mathbf{F}}^{\perp} \mathbf{X}_{\text{new}} \mathbf{P}_{\Theta}^{\perp} + \mathbf{X} \mathbf{P}_{\Theta}$$

- ▶ \mathbf{X}_{new} : fresh i.i.d. $\mathcal{N}(0, 1)$ matrix, **independent** of training data
- ▶ $\mathbf{P}_{\Theta}, \mathbf{P}_{\mathbf{F}}$: projectors onto parameter/gradient subspaces

This technique has been extensively studied in AMP literature⁹.

Result:

$$\mathbf{H}_j(t) \stackrel{d}{=} \mathbf{X}_{\text{new}}^{\top} \mathbf{G}_j(t) \mathbf{X}_{\text{new}} + \text{finite-rank perturbation}$$

⇒ Bulk follows **generalized Marchenko-Pastur law**

⇒ Outliers come from finite-rank perturbation

⁹[Bayati-Montanari '10, Donoho-Montanari '13]

Proof sketch: Step 1 – Gaussian conditioning

Gaussian conditioning decomposition:

$$\mathbf{X} \stackrel{d}{=} \mathbf{P}_{\mathbf{F}}^{\perp} \mathbf{X}_{\text{new}} \mathbf{P}_{\Theta}^{\perp} + \mathbf{X} \mathbf{P}_{\Theta}$$

- ▶ \mathbf{X}_{new} : fresh i.i.d. $\mathcal{N}(0, 1)$ matrix, **independent** of training data
- ▶ \mathbf{P}_{Θ} , $\mathbf{P}_{\mathbf{F}}$: projectors onto parameter/gradient subspaces

This technique has been extensively studied in AMP literature⁹.

Result:

$$\mathbf{H}_j(t) \stackrel{d}{=} \mathbf{X}_{\text{new}}^{\top} \mathbf{G}_j(t) \mathbf{X}_{\text{new}} + \text{finite-rank perturbation}$$

⇒ Bulk follows **generalized Marchenko-Pastur law**

⇒ Outliers come from finite-rank perturbation

⁹[Bayati-Montanari '10, Donoho-Montanari '13]

Proof sketch: Step 1 – Gaussian conditioning

Gaussian conditioning decomposition:

$$\mathbf{X} \stackrel{d}{=} \mathbf{P}_{\mathbf{F}}^{\perp} \mathbf{X}_{\text{new}} \mathbf{P}_{\Theta}^{\perp} + \mathbf{X} \mathbf{P}_{\Theta}$$

- ▶ \mathbf{X}_{new} : fresh i.i.d. $\mathcal{N}(0, 1)$ matrix, **independent** of training data
- ▶ \mathbf{P}_{Θ} , $\mathbf{P}_{\mathbf{F}}$: projectors onto parameter/gradient subspaces

This technique has been extensively studied in AMP literature⁹.

Result:

$$\mathbf{H}_j(t) \stackrel{d}{=} \mathbf{X}_{\text{new}}^{\top} \mathbf{G}_j(t) \mathbf{X}_{\text{new}} + \text{finite-rank perturbation}$$

⇒ Bulk follows **generalized Marchenko-Pastur law**

⇒ Outliers come from finite-rank perturbation

⁹[Bayati-Montanari '10, Donoho-Montanari '13]

Proof sketch: Step 1 – Gaussian conditioning

Gaussian conditioning decomposition:

$$\mathbf{X} \stackrel{d}{=} \mathbf{P}_{\mathbf{F}}^{\perp} \mathbf{X}_{\text{new}} \mathbf{P}_{\Theta}^{\perp} + \mathbf{X} \mathbf{P}_{\Theta}$$

- ▶ \mathbf{X}_{new} : fresh i.i.d. $\mathcal{N}(0, 1)$ matrix, **independent** of training data
- ▶ \mathbf{P}_{Θ} , $\mathbf{P}_{\mathbf{F}}$: projectors onto parameter/gradient subspaces

This technique has been extensively studied in AMP literature⁹.

Result:

$$\mathbf{H}_j(\mathbf{t}) \stackrel{d}{=} \mathbf{X}_{\text{new}}^{\top} \mathbf{G}_j(\mathbf{t}) \mathbf{X}_{\text{new}} + \text{finite-rank perturbation}$$

⇒ Bulk follows **generalized Marchenko-Pastur law**

⇒ Outliers come from finite-rank perturbation

⁹[Bayati-Montanari '10, Donoho-Montanari '13]

Proof sketch: Step 1 – Gaussian conditioning

Gaussian conditioning decomposition:

$$\mathbf{X} \stackrel{d}{=} \mathbf{P}_{\mathbf{F}}^{\perp} \mathbf{X}_{\text{new}} \mathbf{P}_{\Theta}^{\perp} + \mathbf{X} \mathbf{P}_{\Theta}$$

- ▶ \mathbf{X}_{new} : fresh i.i.d. $\mathcal{N}(0, 1)$ matrix, **independent** of training data
- ▶ \mathbf{P}_{Θ} , $\mathbf{P}_{\mathbf{F}}$: projectors onto parameter/gradient subspaces

This technique has been extensively studied in AMP literature⁹.

Result:

$$\mathbf{H}_j(\mathbf{t}) \stackrel{d}{=} \mathbf{X}_{\text{new}}^{\top} \mathbf{G}_j(\mathbf{t}) \mathbf{X}_{\text{new}} + \text{finite-rank perturbation}$$

⇒ Bulk follows **generalized Marchenko-Pastur law**

⇒ Outliers come from finite-rank perturbation

⁹[Bayati-Montanari '10, Donoho-Montanari '13]

Proof sketch: Step 2 – Outlier detection

From Step 1: $\mathbf{H}_j(\mathbf{t}) \stackrel{d}{=} \mathbf{X}_{\text{new}}^\top \mathbf{G}_j \mathbf{X}_{\text{new}} + \text{finite-rank}$

Woodbury identity: For $A + UCV$,

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}$$

Apply to resolvent:

$$(\mathbf{H}_j(\mathbf{t}) - zI)^{-1} = R_0(z) - R_0(z) \cdot [\text{low-rank}] \cdot \mathbf{M}_j(z; \mathbf{t})^{-1} \cdot [\text{low-rank}]^\top \cdot R_0(z)$$

where $R_0(z) = (\mathbf{X}_{\text{new}}^\top \mathbf{G}_j \mathbf{X}_{\text{new}} - zI)^{-1}$ is the **bulk resolvent**

Key: z is an outlier eigenvalue $\Leftrightarrow \det(\mathbf{M}_j(z; \mathbf{t})) = 0$

$\mathbf{M}_j(z; \mathbf{t})$ is **finite-dimensional** (size $\sim k + mt$, independent of n , d)

Proof sketch: Step 2 – Outlier detection

From Step 1: $\mathbf{H}_j(\mathbf{t}) \stackrel{d}{=} \mathbf{X}_{\text{new}}^\top \mathbf{G}_j \mathbf{X}_{\text{new}} + \text{finite-rank}$

Woodbury identity: For $\mathbf{A} + \mathbf{UCV}$,

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1}$$

Apply to resolvent:

$$(\mathbf{H}_j(\mathbf{t}) - z\mathbf{I})^{-1} = \mathbf{R}_0(z) - \mathbf{R}_0(z) \cdot [\text{low-rank}] \cdot \mathbf{M}_j(z; \mathbf{t})^{-1} \cdot [\text{low-rank}]^\top \cdot \mathbf{R}_0(z)$$

where $\mathbf{R}_0(z) = (\mathbf{X}_{\text{new}}^\top \mathbf{G}_j \mathbf{X}_{\text{new}} - z\mathbf{I})^{-1}$ is the **bulk resolvent**

Key: z is an outlier eigenvalue $\Leftrightarrow \det(\mathbf{M}_j(z; \mathbf{t})) = 0$

$\mathbf{M}_j(z; \mathbf{t})$ is **finite-dimensional** (size $\sim k + mt$, independent of n, d)

Proof sketch: Step 2 – Outlier detection

From Step 1: $\mathbf{H}_j(\mathbf{t}) \stackrel{d}{=} \mathbf{X}_{\text{new}}^\top \mathbf{G}_j \mathbf{X}_{\text{new}} + \text{finite-rank}$

Woodbury identity: For $\mathbf{A} + \mathbf{UCV}$,

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{U} (\mathbf{C}^{-1} + \mathbf{VA}^{-1} \mathbf{U})^{-1} \mathbf{VA}^{-1}$$

Apply to resolvent:

$$(\mathbf{H}_j(\mathbf{t}) - z\mathbf{I})^{-1} = \mathbf{R}_0(z) - \mathbf{R}_0(z) \cdot [\text{low-rank}] \cdot \mathbf{M}_j(z; \mathbf{t})^{-1} \cdot [\text{low-rank}]^\top \cdot \mathbf{R}_0(z)$$

where $\mathbf{R}_0(z) = (\mathbf{X}_{\text{new}}^\top \mathbf{G}_j \mathbf{X}_{\text{new}} - z\mathbf{I})^{-1}$ is the **bulk resolvent**

Key: z is an outlier eigenvalue $\Leftrightarrow \det(\mathbf{M}_j(z; \mathbf{t})) = 0$

$\mathbf{M}_j(z; \mathbf{t})$ is **finite-dimensional** (size $\sim k + mt$, independent of n, d)

Proof sketch: Step 2 – Outlier detection

From Step 1: $\mathbf{H}_j(\mathbf{t}) \stackrel{d}{=} \mathbf{X}_{\text{new}}^\top \mathbf{G}_j \mathbf{X}_{\text{new}} + \text{finite-rank}$

Woodbury identity: For $\mathbf{A} + \mathbf{UCV}$,

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1}$$

Apply to resolvent:

$$(\mathbf{H}_j(\mathbf{t}) - z\mathbf{I})^{-1} = \mathbf{R}_0(z) - \mathbf{R}_0(z) \cdot [\text{low-rank}] \cdot \mathbf{M}_j(z; \mathbf{t})^{-1} \cdot [\text{low-rank}]^\top \cdot \mathbf{R}_0(z)$$

where $\mathbf{R}_0(z) = (\mathbf{X}_{\text{new}}^\top \mathbf{G}_j \mathbf{X}_{\text{new}} - z\mathbf{I})^{-1}$ is the **bulk resolvent**

Key: z is an outlier eigenvalue $\Leftrightarrow \det(\mathbf{M}_j(z; \mathbf{t})) = 0$

$\mathbf{M}_j(z; \mathbf{t})$ is **finite-dimensional** (size $\sim k + mt$, independent of n, d)

Proof sketch: Step 2 – Block structure

Concentration: $\mathbf{M}_j(z; t) \xrightarrow{p} M_j^\infty(z; t)$ (deterministic limit)

Fortunately: $M_j^\infty(z; t)$ has **block-diagonal structure!**

$$\det(M_j^\infty(z; t)) = 0 \quad \Leftrightarrow \quad \det(M_{j,H}^\infty) = 0 \text{ or } \det(M_{j,R}^\infty) = 0$$

- ▶ $M_{j,H}^\infty$: Hard block (size $r \times r$, independent of t)
 $\Rightarrow \det(M_{j,H}^\infty) = 0$ gives our outlier equation
- ▶ $M_{j,R}^\infty$: Rest block
 \Rightarrow Can prove: solutions have **no informative eigenvectors**

Proof sketch: Step 2 – Block structure

Concentration: $\mathbf{M}_j(z; t) \xrightarrow{P} \mathbf{M}_j^\infty(z; t)$ (deterministic limit)

Fortunately: $\mathbf{M}_j^\infty(z; t)$ has **block-diagonal structure!**

$$\det(\mathbf{M}_j^\infty(z; t)) = 0 \quad \Leftrightarrow \quad \det(\mathbf{M}_{j,H}^\infty) = 0 \text{ or } \det(\mathbf{M}_{j,R}^\infty) = 0$$

- ▶ $\mathbf{M}_{j,H}^\infty$: Hard block (size $r \times r$, independent of t)
 $\Rightarrow \det(\mathbf{M}_{j,H}^\infty) = 0$ gives our outlier equation
- ▶ $\mathbf{M}_{j,R}^\infty$: Rest block
 \Rightarrow Can prove: solutions have no informative eigenvectors

Proof sketch: Step 2 – Block structure

Concentration: $M_j(z; t) \xrightarrow{P} M_j^\infty(z; t)$ (deterministic limit)

Fortunately: $M_j^\infty(z; t)$ has **block-diagonal structure!**

$$\det(M_j^\infty(z; t)) = 0 \quad \Leftrightarrow \quad \det(M_{j,H}^\infty) = 0 \text{ or } \det(M_{j,R}^\infty) = 0$$

- ▶ $M_{j,H}^\infty$: **Hard block** (size $r \times r$, independent of t)
 $\Rightarrow \det(M_{j,H}^\infty) = 0$ gives our **outlier equation**
- ▶ $M_{j,R}^\infty$: **Rest block**
 \Rightarrow Can prove: solutions have **no informative eigenvectors**

Proof sketch: Step 2 – Block structure

Concentration: $M_j(z; t) \xrightarrow{P} M_j^\infty(z; t)$ (deterministic limit)

Fortunately: $M_j^\infty(z; t)$ has **block-diagonal structure!**

$$\det(M_j^\infty(z; t)) = 0 \quad \Leftrightarrow \quad \det(M_{j,H}^\infty) = 0 \text{ or } \det(M_{j,R}^\infty) = 0$$

- ▶ $M_{j,H}^\infty$: **Hard block** (size $r \times r$, independent of t)
 $\Rightarrow \det(M_{j,H}^\infty) = 0$ gives our **outlier equation**
- ▶ $M_{j,R}^\infty$: **Rest block**
 \Rightarrow Can prove: solutions have **no informative eigenvectors**

Proof sketch: Step 3 – Eigenvector alignment

Residue formula:

$$\xi \xi^\top = -\frac{1}{2\pi i} \oint_{\gamma} (\mathbf{H}_j(t) - zI)^{-1} dz$$

Correlation:

$$\|\Theta_{*H}^\top \xi\|^2 = -\frac{1}{2\pi i} \oint_{\gamma} \text{Tr}(\Theta_{*H}^\top (\mathbf{H}_j(t) - zI)^{-1} \Theta_{*H}) dz$$

Simplification:

- ▶ Substitute resolvent expansion from the last step
- ▶ $R_0(z)$ analytic inside $\gamma \Rightarrow$ only pole from $M_j(z; t)^{-1}$
- ▶ Replace $M_j(z; t)$ by limit $M_j^\infty(z; t)$, apply residue theorem
- ▶ Use block structure of $M_j^\infty \Rightarrow$ final formula

Proof sketch: Step 3 – Eigenvector alignment

Residue formula:

$$\xi \xi^\top = -\frac{1}{2\pi i} \oint_{\gamma} (\mathbf{H}_j(t) - zI)^{-1} dz$$

Correlation:

$$\|\Theta_{*H}^\top \xi\|^2 = -\frac{1}{2\pi i} \oint_{\gamma} \text{Tr}(\Theta_{*H}^\top (\mathbf{H}_j(t) - zI)^{-1} \Theta_{*H}) dz$$

Simplification:

- ▶ Substitute resolvent expansion from the last step
- ▶ $R_0(z)$ analytic inside $\gamma \Rightarrow$ only pole from $M_j(z; t)^{-1}$
- ▶ Replace $M_j(z; t)$ by limit $M_j^\infty(z; t)$, apply residue theorem
- ▶ Use block structure of $M_j^\infty \Rightarrow$ final formula

Proof sketch: Step 3 – Eigenvector alignment

Residue formula:

$$\xi \xi^\top = -\frac{1}{2\pi i} \oint_{\gamma} (\mathbf{H}_j(\mathbf{t}) - zI)^{-1} dz$$

Correlation:

$$\|\Theta_{*H}^\top \xi\|^2 = -\frac{1}{2\pi i} \oint_{\gamma} \text{Tr}(\Theta_{*H}^\top (\mathbf{H}_j(\mathbf{t}) - zI)^{-1} \Theta_{*H}) dz$$

Simplification:

- ▶ Substitute resolvent expansion from the last step
- ▶ $R_0(z)$ analytic inside $\gamma \Rightarrow$ only pole from $\mathbf{M}_j(z; \mathbf{t})^{-1}$
- ▶ Replace $\mathbf{M}_j(z; \mathbf{t})$ by limit $\mathbf{M}_j^\infty(z; \mathbf{t})$, apply residue theorem
- ▶ Use block structure of $\mathbf{M}_j^\infty \Rightarrow$ final formula

Proof sketch: Step 3 – Eigenvector alignment

Residue formula:

$$\xi \xi^\top = -\frac{1}{2\pi i} \oint_{\gamma} (\mathbf{H}_j(t) - zI)^{-1} dz$$

Correlation:

$$\|\Theta_{*H}^\top \xi\|^2 = -\frac{1}{2\pi i} \oint_{\gamma} \text{Tr}(\Theta_{*H}^\top (\mathbf{H}_j(t) - zI)^{-1} \Theta_{*H}) dz$$

Simplification:

- ▶ Substitute resolvent expansion from the last step
- ▶ $R_0(z)$ analytic inside $\gamma \Rightarrow$ **only pole from $\mathbf{M}_j(z; t)^{-1}$**
- ▶ Replace $\mathbf{M}_j(z; t)$ by limit $\mathbf{M}_j^\infty(z; t)$, apply residue theorem
- ▶ Use block structure of $\mathbf{M}_j^\infty \Rightarrow$ final formula

Proof sketch: Step 3 – Eigenvector alignment

Residue formula:

$$\xi \xi^\top = -\frac{1}{2\pi i} \oint_{\gamma} (\mathbf{H}_j(\mathbf{t}) - zI)^{-1} dz$$

Correlation:

$$\|\Theta_{*H}^\top \xi\|^2 = -\frac{1}{2\pi i} \oint_{\gamma} \text{Tr}(\Theta_{*H}^\top (\mathbf{H}_j(\mathbf{t}) - zI)^{-1} \Theta_{*H}) dz$$

Simplification:

- ▶ Substitute resolvent expansion from the last step
- ▶ $R_0(z)$ analytic inside $\gamma \Rightarrow$ **only pole from $\mathbf{M}_j(z; \mathbf{t})^{-1}$**
- ▶ Replace $\mathbf{M}_j(z; \mathbf{t})$ by limit $M_j^\infty(z; \mathbf{t})$, apply residue theorem
- ▶ Use block structure of $M_j^\infty \Rightarrow$ final formula

Proof sketch: Step 3 – Eigenvector alignment

Residue formula:

$$\xi \xi^\top = -\frac{1}{2\pi i} \oint_{\gamma} (\mathbf{H}_j(t) - zI)^{-1} dz$$

Correlation:

$$\|\Theta_{*H}^\top \xi\|^2 = -\frac{1}{2\pi i} \oint_{\gamma} \text{Tr}(\Theta_{*H}^\top (\mathbf{H}_j(t) - zI)^{-1} \Theta_{*H}) dz$$

Simplification:

- ▶ Substitute resolvent expansion from the last step
- ▶ $R_0(z)$ analytic inside $\gamma \Rightarrow$ **only pole from $\mathbf{M}_j(z; t)^{-1}$**
- ▶ Replace $\mathbf{M}_j(z; t)$ by limit $M_j^\infty(z; t)$, apply residue theorem
- ▶ Use block structure of $M_j^\infty \Rightarrow$ final formula

Grokking explanation

Setup: (We assume for now) only hard directions to learn, $\delta > \delta_{\text{NN}}$

Stage 1: $t = O(1)$

- ▶ Cannot learn hard directions (no linear signal)
- ▶ But $n/d = \delta$ is finite \Rightarrow **overfits a bit training data**
- ▶ Train loss \ll Test loss (generalization gap)

Stage 2: Saddle escape happens

- ▶ Since $\delta > \delta_{\text{NN}}$: Hessian **develops informative outlier**
- ▶ Eventually **learns hard directions**
- ▶ Test loss drops \Rightarrow **Predictable grokking!**

No grokking when $\delta \gg \delta_{\text{NN}}$:

- ▶ No overfitting in Stage 1
- ▶ No generalization gap \Rightarrow no delayed generalization

Grokking explanation

Setup: (We assume for now) only hard directions to learn, $\delta > \delta_{\text{NN}}$

Stage 1: $t = O(1)$

- ▶ Cannot learn hard directions (no linear signal)
- ▶ But $n/d = \delta$ is finite \Rightarrow **overfits a bit training data**
- ▶ Train loss \ll Test loss (generalization gap)

Stage 2: Saddle escape happens

- ▶ Since $\delta > \delta_{\text{NN}}$: Hessian **develops informative outlier**
- ▶ Eventually **learns hard directions**
- ▶ Test loss drops \Rightarrow **Predictable grokking!**

No grokking when $\delta \gg \delta_{\text{NN}}$:

- ▶ No overfitting in Stage 1
- ▶ No generalization gap \Rightarrow no delayed generalization

Grokking explanation

Setup: (We assume for now) only hard directions to learn, $\delta > \delta_{\text{NN}}$

Stage 1: $t = O(1)$

- ▶ Cannot learn hard directions (no linear signal)
- ▶ But $n/d = \delta$ is finite \Rightarrow **overfits a bit training data**
- ▶ Train loss \ll Test loss (generalization gap)

Stage 2: Saddle escape happens

- ▶ Since $\delta > \delta_{\text{NN}}$: Hessian **develops informative outlier**
- ▶ Eventually **learns hard directions**
- ▶ Test loss drops \Rightarrow **Predictable grokking!**

No grokking when $\delta \gg \delta_{\text{NN}}$:

- ▶ No overfitting in Stage 1
- ▶ No generalization gap \Rightarrow no delayed generalization

Grokking explanation

Setup: (We assume for now) only hard directions to learn, $\delta > \delta_{\text{NN}}$

Stage 1: $t = O(1)$

- ▶ Cannot learn hard directions (no linear signal)
- ▶ But $n/d = \delta$ is finite \Rightarrow **overfits a bit training data**
- ▶ Train loss \ll Test loss (generalization gap)

Stage 2: Saddle escape happens

- ▶ Since $\delta > \delta_{\text{NN}}$: Hessian **develops informative outlier**
- ▶ Eventually **learns hard directions**
- ▶ Test loss drops \Rightarrow **Predictable grokking!**

No grokking when $\delta \gg \delta_{\text{NN}}$:

- ▶ No overfitting in Stage 1
- ▶ No generalization gap \Rightarrow no delayed generalization

Discussion

Conclusion

Main contributions:

- **Precise threshold δ_{NN}** for feature learning (the first hard direction) via gradient descent on two-layer NNs
 - Depends on target, loss, activation, width, initialization
 - **Universally sharp** across broad problem setups
- **Spectral mechanism:** Learning hard directions \Leftrightarrow negative outliers of Hessian aligned with hard subspace
- $\delta_{\text{NN}} > \delta_{\text{alg}}$ in general; gap depends on architecture and algorithm
- **Grokking explanation:** Delayed generalization predictably occurs when and only when δ is moderately above δ_{NN}

Conclusion

Main contributions:

- 1 **Precise threshold δ_{NN}** for feature learning (the first hard direction) via gradient descent on two-layer NNs
 - ▶ Depends on target, loss, activation, width, initialization
 - ▶ **Efficiently computable** from problem setup
- 2 **Spectral mechanism:** Learning hard directions \Leftrightarrow negative outliers of Hessian aligned with hard subspace
- 3 $\delta_{\text{NN}} > \delta_{\text{alg}}$ in general; gap depends on architecture and algorithm
- 4 **Grokking explanation:** Delayed generalization predictably occurs when and only when δ is moderately above δ_{NN}

Conclusion

Main contributions:

- 1 **Precise threshold δ_{NN}** for feature learning (the first hard direction) via gradient descent on two-layer NNs
 - ▶ Depends on target, loss, activation, width, initialization
 - ▶ **Efficiently computable** from problem setup
- 2 **Spectral mechanism:** Learning hard directions \Leftrightarrow negative outliers of Hessian aligned with hard subspace
- 3 $\delta_{\text{NN}} > \delta_{\text{alg}}$ in general; gap depends on architecture and algorithm
- 4 **Grokking explanation:** Delayed generalization predictably occurs when and only when δ is moderately above δ_{NN}

Conclusion

Main contributions:

- 1 **Precise threshold δ_{NN}** for feature learning (the first hard direction) via gradient descent on two-layer NNs
 - ▶ Depends on target, loss, activation, width, initialization
 - ▶ **Efficiently computable** from problem setup
- 2 **Spectral mechanism:** Learning hard directions \Leftrightarrow negative outliers of Hessian aligned with hard subspace
- 3 $\delta_{\text{NN}} > \delta_{\text{alg}}$ in general; gap depends on architecture and algorithm
- 4 **Grokking explanation:** Delayed generalization predictably occurs when and only when δ is moderately above δ_{NN}

Conclusion

Main contributions:

- 1 **Precise threshold** δ_{NN} for feature learning (the first hard direction) via gradient descent on two-layer NNs
 - ▶ Depends on target, loss, activation, width, initialization
 - ▶ **Efficiently computable** from problem setup
- 2 **Spectral mechanism:** Learning hard directions \Leftrightarrow negative outliers of Hessian aligned with hard subspace
- 3 $\delta_{\text{NN}} > \delta_{\text{alg}}$ in general; gap depends on architecture and algorithm
- 4 **Grokking explanation:** Delayed generalization predictably occurs when and only when δ is moderately above δ_{NN}

Conclusion

Main contributions:

- ① **Precise threshold δ_{NN}** for feature learning (the first hard direction) via gradient descent on two-layer NNs
 - ▶ Depends on target, loss, activation, width, initialization
 - ▶ **Efficiently computable** from problem setup
- ② **Spectral mechanism:** Learning hard directions \Leftrightarrow negative outliers of Hessian aligned with hard subspace
- ③ $\delta_{\text{NN}} > \delta_{\text{alg}}$ in general; gap depends on architecture and algorithm
- ④ **Grokking explanation:** Delayed generalization predictably occurs when and only when δ is moderately above δ_{NN}

Conclusion

Main contributions:

- ① **Precise threshold δ_{NN}** for feature learning (the first hard direction) via gradient descent on two-layer NNs
 - ▶ Depends on target, loss, activation, width, initialization
 - ▶ **Efficiently computable** from problem setup
- ② **Spectral mechanism:** Learning hard directions \Leftrightarrow negative outliers of Hessian aligned with hard subspace
- ③ $\delta_{\text{NN}} > \delta_{\text{alg}}$ in general; gap depends on architecture and algorithm
- ④ **Grokking explanation:** Delayed generalization predictably occurs **when and only when δ is moderately above δ_{NN}**

Future directions

- ① Results for finite $m = O(1)$
- ② Training both layers simultaneously
- ③ Beyond $O(1)$ time: Precise dynamics for learning hard directions (requires diverging time)
- ④ Thresholds for learning the whole subspace
- ⑤ Non-isotropic covariates?
- ⑥ Hierarchical feature learning in deeper models?

Thanks for listening!

Questions?

[arXiv:2602.01434](https://arxiv.org/abs/2602.01434)