

Understanding modern machine learning models through the lens of high-dimensional statistics

Denny Wu

Center for Data Science, New York University

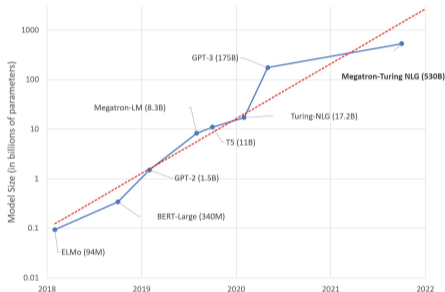
Center for Computational Mathematics, Flatiron Institute



High-dimensionality of Modern ML Systems

Modern ML tasks require searching over a *high-dimensional* parameter space.

Curse of dimensionality? Larger neural networks often achieve better performance.



LLM parameter count (Hugging Face blogpost)

Overparameterization: # parameters > #training data.

High-dimensionality of Modern ML Systems

Modern ML tasks require searching over a *high-dimensional* parameter space.

Curse of dimensionality? Larger neural networks often achieve better performance.

Understanding the success of deep learning

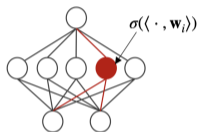
- (i) **Optimization:** standard gradient-based methods work, despite the non-convexity.
 - benefit of overparameterization (NTK, mean-field, etc.)
- (ii) **Generalization:** model generalizes well, despite the overparameterization.
 - implicit regularization, benign overfitting.
- (iii) **Why neural networks?** NN often outperforms classical methods (e.g., kernels).
 - adaptivity, representation (feature) learning.

My research: quantitative understanding of (i)-(iii) via *high-dimensional statistics*.

Mathematical Models for High-dimensional Problems

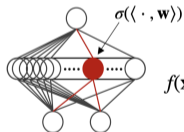
Intuition: theoretical analysis may simplify if we take the *dimensionality to infinity*.

Scaling ① – Large Width Limit



$$f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \sigma(\langle \mathbf{x}, \mathbf{w}_i \rangle)$$

$$N \rightarrow \infty$$



$$f(\mathbf{x}) = \mathbb{E}_p[\sigma(\langle \mathbf{x}, \mathbf{w} \rangle)] = \int \sigma(\langle \mathbf{x}, \mathbf{w} \rangle) p(\mathbf{w}) d\mathbf{w}$$

For convex loss L , learning is

☹ **non-convex** w.r.t. \mathbf{w}_i

☺ **convex** w.r.t. distribution p

Perspective: study optimization in the space of measures

(Wasserstein gradient flow, functional inequalities (LSI), etc.)

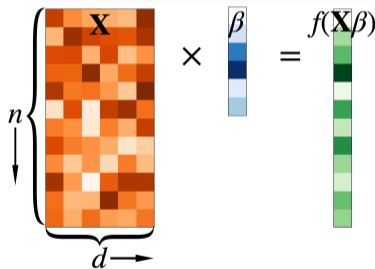
- Convergence rate of mean-field Langevin dynamics and propagation of chaos [NWS22][SWN23]
- Learnability guarantees for low-dimensional target functions [SWO+23][NOS+23]
- New algorithms for optimization in the space of measures [NWS21] [OSN+22] [NOW+23]
-

Mathematical Models for High-dimensional Problems

Intuition: theoretical analysis may simplify if we take the *dimensionality to infinity*.

Scaling ② – Proportional Asymptotic Limit

$n, d \rightarrow \infty, d/n = \gamma \in (0, \infty)$



Diverging dimensionality & *fixed aspect ratio*.

- Captures the *overparameterized* regime (by setting $\gamma > 1$)

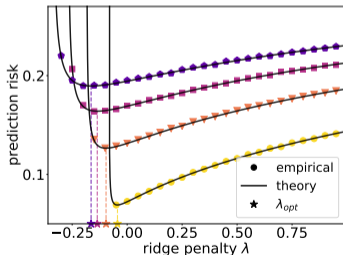
Performance of simple ML models can be **precisely analyzed** via *random matrix theory* (the study of large-dimensional matrices with certain *random* structures)

This talk: two examples of precise analysis using random matrix theory (RMT).
(i) optimal regularization in linear regression. (ii) feature learning in neural network.

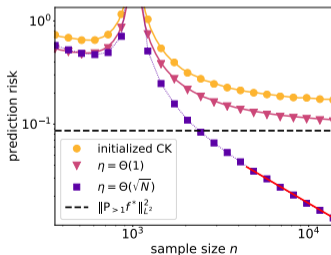
Precise Analysis of Learning in High Dimensions

What are the advantages of a **precise analysis**?

- Enables *accurate comparison* between estimators/algorithms.
 - positive vs. negative ridge penalty, gradient descent vs. natural gradient, etc.
- Captures refined properties of the learning curve.
 - phase transitions, (non-)monotonicity, etc.



"Negative ridge" phenomenon.



Benefit of representation learning.

Ridge Regression in High Dimensions

Problem Setting & Assumptions

- **Data Generation:** $y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_* + \varepsilon_i$, $1 \leq i \leq n$. $\mathbf{x}_i \in \mathbb{R}^d$.
i.i.d. label noise satisfies $\mathbb{E}[\varepsilon] = 0$, $\text{Var}(\varepsilon) = \sigma_\varepsilon^2$.
- **Random Design:** $\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \boldsymbol{\Sigma})$. Also holds for \mathbf{x}_i with bounded $(4+\epsilon)$ moment
- **Signal (Ground Truth):** $\boldsymbol{\beta}_*$ can be both *fixed* or *random* (i.e., $\mathbb{E}[\boldsymbol{\beta}_* \boldsymbol{\beta}_*^\top] = \boldsymbol{\Sigma}_\beta$)
- **Proportional Asymptotics:** $n, d \rightarrow \infty$, $d/n \rightarrow \gamma \in (0, \infty)$.

Ridge regression estimator: $\hat{\boldsymbol{\beta}}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^\dagger \mathbf{X}^\top \mathbf{y}$.

- **Goal:** compute the *prediction risk* (test error) $\mathcal{R}(\lambda) = \mathbb{E}(y - \mathbf{x}^\top \hat{\boldsymbol{\beta}}_\lambda)^2$.

Remark: When $\lambda \geq 0$, $\hat{\boldsymbol{\beta}}_\lambda = \arg \min_{\boldsymbol{\beta}} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 + \lambda \|\boldsymbol{\beta}\|^2$ (Gaussian prior)

Theorem ([WX20] Precise generalization error of ridge regression)

The test error admits a bias-variance decomposition $\mathcal{R}(\lambda) = \mathcal{B}(\lambda) + \mathcal{V}(\lambda)$, where

$$\mathcal{B}(\lambda) \xrightarrow{p} \frac{\partial \kappa_\lambda}{\partial \lambda} \cdot \kappa_\lambda^2 \langle \beta_*, \Sigma(\Sigma + \kappa_\lambda \mathbf{I})^{-2} \beta_* \rangle, \quad \mathcal{V}(\lambda) \xrightarrow{p} \sigma_\varepsilon^2 \frac{\partial \kappa_\lambda}{\partial \lambda},$$

and $\kappa_\lambda \geq \lambda$ is the **effective regularization** given by the non-negative solution of

$$\frac{1}{n} \text{Tr}(\Sigma(\Sigma + \kappa_\lambda \mathbf{I})^{-1}) = 1 - \frac{\lambda}{\kappa_\lambda}.$$

- **Bias** $\mathcal{B}(\lambda)$: learning of *signal* β_* .
- **Variance** $\mathcal{V}(\lambda)$: "overfitting" to *label noise*.

Given *eigendecomposition* $\Sigma = \sum_i \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$, prediction risk $\mathcal{R}(\lambda)$ depends on:

- **Capacity condition:** eigenvalues of the population covariance $\{\lambda_i\}_{i=1}^d$.
- **Source condition:** projection of signal (teacher) β_* onto the feature eigenbasis $\{\rho_i\}_{i=1}^d$, where $\rho_i = \langle \beta_*, \mathbf{u}_i \rangle$.

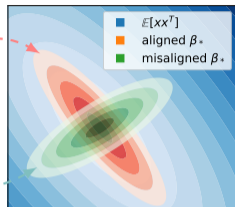
Alignment between Features and Signal

- Aligned feature & signal: large $\lambda_i \Leftrightarrow$ large $\langle \beta_*, \mathbf{u}_i \rangle$

☺ Features are well-engineered \Rightarrow easy problem

- Misaligned feature & signal: large $\lambda_i \Leftrightarrow$ small $\langle \beta_*, \mathbf{u}_i \rangle$

☹ Features are uninformative \Rightarrow hard problem



Theorem ([WX20] Sign of Optimal Ridge Penalty λ_{opt})

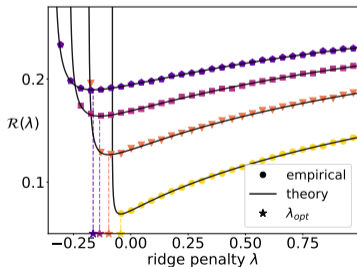
Recall that $\gamma = d/n$.

- ☐ $\gamma < 1$ (**underparameterized**): $\lambda_{\text{opt}} \geq 0$ in all cases.
- ☐ $\gamma > 1$ (**overparameterized**): the sign of λ_{opt} depends on the *alignment* between the features and the signal.

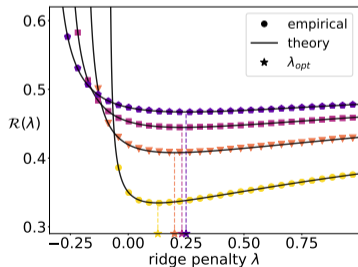
The “Negative Ridge” Phenomenon

Corollary ([WX20] Sign of λ_{opt} in the Overparameterized Regime)

- ❑ **Negative** λ is beneficial under alignment (informative features); hence interpolation ($\lambda = 0$) can be optimal even if $\sigma_\epsilon > 0$.
- ❑ **Positive** λ is beneficial under misalignment (hard problem), even in the absence of label noise ($\sigma_\epsilon = 0$).



Aligned: $\lambda_{opt} < 0$.



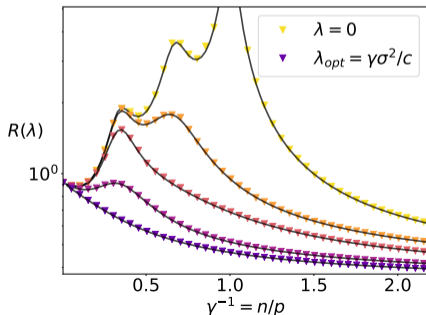
Misaligned: $\lambda_{opt} > 0$.

Regularization Suppresses “Multiple Descent”

Without appropriate regularization, $\mathcal{R}(\lambda)$ may exhibit *multiple peaks*...

Theorem ([WX20] Monotonicity of $\mathcal{R}(\lambda_{opt})$)

Given $\mathbb{E}[\beta_*\beta_*^\top] \propto \mathbf{I}$ (isotropic prior), the optimally regularized prediction risk $\mathcal{R}(\lambda_{opt})$ is a **decreasing** function of $\gamma^{-1} = n/d \in (0, \infty)$.



Message: if we tune λ , more training data *always helps* the test performance.

Implication I: Implicit Bias of Optimizers

$$\text{Update rule: } \theta_{t+1} = \theta_t - \eta \mathbf{P}(\theta_t) \nabla_{\theta_t} L(\theta_t), \quad t = 0, 1, \dots$$

Geometric Intuition: \mathbf{P} alleviates pathological curvature and speed up optimization.

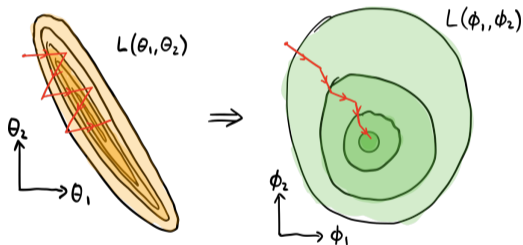


Figure from Xanadu blog post.

Question: in the *interpolation setting* (i.e. absence of explicit regularization), how does preconditioning influence the **generalization** performance?

Implicit Bias in Overparameterized Linear Regression

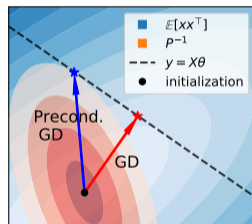
Theoretical Setting: preconditioned gradient descent (flow) on the *overparameterized* least squares objective: $L(\beta) = \frac{1}{n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$.

Implicit Bias ($t \rightarrow \infty$):

- **Gradient descent:** min ℓ_2 norm solution.
- **Preconditioned GD:** for *time-independent* and full-rank \mathbf{P} , min $\|\beta\|_{\mathbf{P}^{-1}}$ norm solution.

Example.

Natural gradient descent with *population Fisher*: $\mathbf{P} = \Sigma^{-1}$



- **Goal I:** use the asymptotic risk formulae (taking $\lambda \rightarrow 0$) to *precisely* compare the generalization of GD vs. NGD.
- **Goal II:** validate our predictions in neural network experiments.

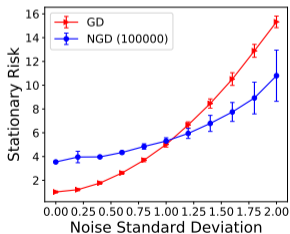
Comparison of Generalization Performance

Theorem ([ABG+21] Prediction Risk of GD vs. NGD)

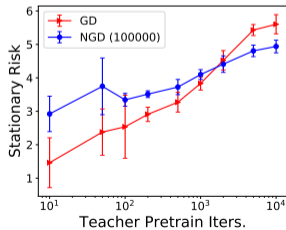
- **Variance** $\lim_{\lambda \rightarrow 0} \mathcal{V}(\lambda)$: **NGD** (population, $\mathbf{P} = \Sigma^{-1}$) is optimal.
- **Bias** $\lim_{\lambda \rightarrow 0} \mathcal{B}(\lambda)$: **GD** generalizes better when signal is isotropic ($\Sigma_{\beta} = \mathbf{I}$); **NGD** generalizes better under misalignment (“difficult problem”).

Remark: bias-variance tradeoff achieved by “interpolating” between optimizers.

Two-layer MLP: student-teacher setup (CIFAR-10)



Label noise (variance).



Misspecification (bias).

Implication II: Beyond Gaussian Features

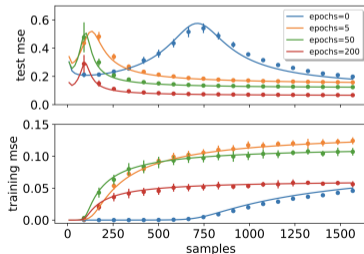
Question: does our risk formula have predictive power in *practical settings*, e.g., neural network representations?

- **Decomposition of kernel:** $k(\mathbf{x}, \mathbf{y}) = \sum_i \lambda_i \phi_i(\mathbf{x}) \phi_i(\mathbf{y})$.
- **Decomposition of target function:** $f^*(\mathbf{x}) = \sum_i \rho_i \phi_i(\mathbf{x})$, $\rho_i = \langle \phi_i, f^* \rangle_{L^2}$.
- **Leap of faith:** estimate $\{\lambda_i, \rho_i\}_{i=1}^{\infty}$ from data, and plug in the risk formulae.

Universality: RMT prediction empirically accurate for many feature maps, including **trained neural network features**.

Observation: trained NN achieves lower risk
⇒ advantage of *representation learning*.

❑ **Spoiler:** this benefit will be precisely analyzed!



[Loureiro et al. 2021]

Two-layer Neural Network

$$f_{\text{NN}}(\mathbf{x}) = \frac{1}{\sqrt{N}} \sum_{i=1}^N a_i \sigma(\mathbf{x}^\top \mathbf{w}_i) = \frac{1}{\sqrt{N}} \mathbf{a}^\top \sigma(\mathbf{W}^\top \mathbf{x}).$$

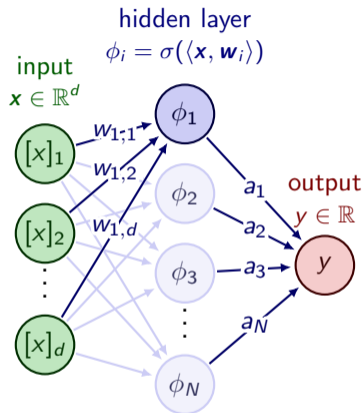
- Trainable parameters: $\mathbf{W} \in \mathbb{R}^{d \times N}$, $\mathbf{a} \in \mathbb{R}^N$.
- Element-wise nonlinearity: $\sigma : \mathbb{R} \rightarrow \mathbb{R}$.

Proportional asymptotic limit:

$$n, d, N \rightarrow \infty, n/d \rightarrow \psi_1, N/d \rightarrow \psi_2, \\ \text{where } \psi_1, \psi_2 \in (0, \infty).$$

- Increase $\psi_1 \Rightarrow$ larger **sample size**.
- Increase $\psi_2 \Rightarrow$ **overparameterization**.

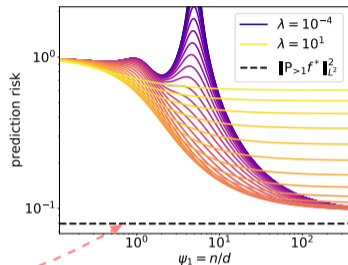
Motivation: rigorously show that the *learned representation* (via gradient descent) achieves better performance in the proportional limit.



Prior Works: Asymptotics of Random Features Model

Fix 1st layer W at initialization, learn 2nd layer a
⇒ **random features (RF)** model.

- ☺ Prediction risk precisely characterized in the *proportional regime* via **random matrix theory**.
- ☹ The (*nonlinear*) RF estimator cannot even outperform **linear functions** on the input...



$P_{>1}$ is the projector orthogonal to constant & linear functions in $L_2(P_x)$

Where does this gap come from?

Feature (representation) learning!

☐ When W is optimized, NN can “adapt” to data and learn useful features.

- Mei and Montanari, 2019. *The generalization error of random features regression: Precise asymptotics and double descent curve.*
- Gerace et al., 2020. *Generalisation error in learning with random features and the hidden manifold model.*

Feature Learning via One Gradient Descent Step

“Early Phase” Feature Learning: Does the first gradient descent step on the first-layer \mathbf{W} already learn useful representations?

- **One-step GD on 1st Layer.** Gradient update $\mathbf{W}_1 = \mathbf{W}_0 + \eta\sqrt{N} \cdot \mathbf{G}$, where

$$\mathbf{G} = -\nabla_{\mathbf{W}} \left[\frac{1}{n} \sum_{i=1}^n (y_i - f_{\text{NN}}(\mathbf{x}_i))^2 \right] = \frac{1}{n} \mathbf{X}^\top \left[\left(\frac{1}{\sqrt{N}} (y - f_{\text{NN}}(\mathbf{X})) \mathbf{a}^\top \right) \odot \sigma'(\mathbf{X} \mathbf{W}_0) \right].$$

- **Ridge Regression for 2nd Layer.** Regression using *trained* kernel features:

$$\hat{\mathbf{a}}_\lambda = \operatorname{argmin}_{\mathbf{a}} \left\{ \frac{1}{n} \|\tilde{\mathbf{y}} - \Phi \mathbf{a}\|^2 + \frac{\lambda}{N} \|\mathbf{a}\|^2 \right\}, \quad \Phi := \frac{1}{\sqrt{N}} \sigma(\tilde{\mathbf{X}} \mathbf{W}_1) \in \mathbb{R}^{n \times N}.$$

Denote $f_{\text{GD}}^\lambda(\mathbf{x}) = \frac{1}{\sqrt{N}} \hat{\mathbf{a}}_\lambda^\top \sigma(\mathbf{W}_1^\top \mathbf{x})$, prediction risk: $\mathcal{R}_{\text{GD}}(\lambda) = \mathcal{R}(f_{\text{GD}}^\lambda)$.

Goal: Precise analysis of $\mathcal{R}_{\text{GD}}(\lambda)$ to show its improvement over the initialized RF $\mathcal{R}_{\text{RF}}(\lambda)$, and potentially over the *kernel lower bound* $\|\mathbb{P}_{>1} f^*\|_{L^2}^2$.

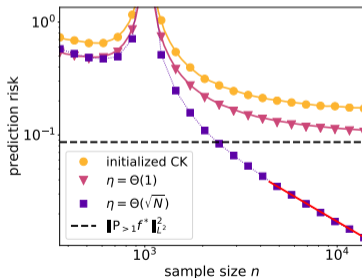
Our Results: Precise Asymptotics of Feature Learning

- **Student-teacher Setup.** $y_i = f^*(\langle \mathbf{x}_i, \boldsymbol{\beta}_* \rangle) + \varepsilon_i$, where $\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I})$.
- **Gaussian Initialization.** $[\mathbf{W}_0]_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1/d)$, $[\mathbf{a}]_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1/N)$.

Small lr ("lazy" regime): $\eta = \Theta(1) \Rightarrow |[\mathbf{W}_1 - \mathbf{W}_0]_{ij}| \ll |[\mathbf{W}_0]_{ij}|$

Large lr (μP scaling): $\eta = \Theta(\sqrt{N}) \Rightarrow |[\mathbf{W}_1 - \mathbf{W}_0]_{ij}| \asymp |[\mathbf{W}_0]_{ij}|$

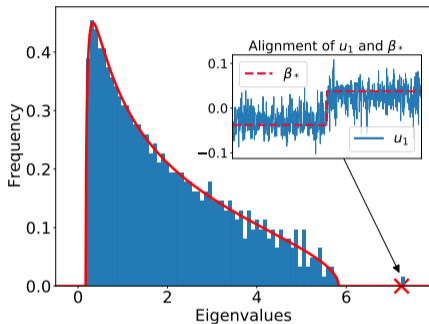
- **Small lr $\eta = \Theta(1)$** : trained kernel **always improve** upon the initial RF estimator, but the model remains "*linear*".
- **Large lr $\eta = \Theta(\sqrt{N})$** : regression on trained features can learn certain **nonlinear** f^* .



- Jacot et al, 2018. *Neural tangent kernel: convergence and generalization in neural networks.*
- Yang and Hu, 2021. *Feature learning in infinite-width neural networks.*

A Spiked Model for the Trained Weight Matrix

Challenge: learned W_1 *no longer i.i.d.*; can we still apply RMT tools?



Blue: empirical simulation

Red: analytic prediction

(*BBP Phase Transition*)

- $\sigma = \tanh$, $f^*(x) = \text{ReLU}(\langle x, \beta_* \rangle)$.
- Teacher $\beta_* \propto [-1_{d/2}; 1_{d/2}]$.

Observation: after *one feature learning step* on the first-layer W :

- The **bulk** of the spectrum of W_1 remains unchanged
- A **spike** (x) appears in W_1 , which aligns with signal β^*

$\eta = \Theta(1)$ – Precise Analysis via Gaussian Equivalence

Intuition – Universality: replace *nonlinear* NN features with *linear* Gaussian features with *matching first two moments* does not change the risk.

- NN (nonlinear): $\phi_{\text{NN}}(\mathbf{x}) = \frac{1}{\sqrt{N}}\sigma(\mathbf{W}^\top \mathbf{x})$.
- GE (linear): $\phi_{\text{GE}}(\mathbf{x}) = \frac{1}{\sqrt{N}}\left(\mu_1 \mathbf{W}^\top \mathbf{x} + \mu_2 \mathbf{z}\right)$, $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$.

$$\text{where } \mu_1 = \mathbb{E}[z\sigma(z)], \mu_2 = \sqrt{\mathbb{E}[\sigma(z)^2] - \mu_1^2} \Rightarrow \mathbb{E}[\phi_{\text{NN}}\phi_{\text{NN}}^\top] = \mathbb{E}[\phi_{\text{GE}}\phi_{\text{GE}}^\top]$$

Theorem ([BES+22] Gaussian Equivalence for Trained Features)

After one feature learning step on \mathbf{W} with small learning rate $\eta = \Theta(1)$,

$$|\mathcal{R}_{\text{GD}}(\lambda) - \mathcal{R}_{\text{GE}}(\lambda)| = o_{d,\mathbb{P}}(1), \text{ for } \lambda > 0.$$

Implications of Gaussian Equivalence (GET):

- We may equivalently compute \mathcal{R}_{GE} , which can be handled via RMT tools ☺
- The nonlinear NN model achieves the same performance as a linear model ☹

Precise Characterization of Feature Learning

Theorem ([BES+22] Benefit of Feature Learning)

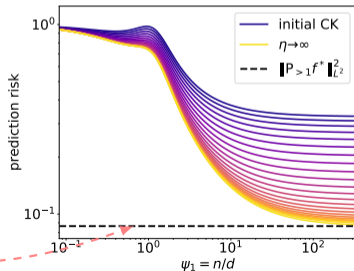
The prediction risk difference $\delta := \lim_{n,d,N \rightarrow \infty} \mathcal{R}_{\text{RF}}(\lambda) - \mathcal{R}_{\text{GD}}(\lambda)$ is

- a **non-negative** function of $\eta, \lambda, \psi_1, \psi_2 \in (0, +\infty)$;
- an **increasing** function with respect to learning rate η .

Provable improvement over the initial RF model!

Observations:

- For $\eta = \Theta(1)$, feature learning always helps.
- Larger step size \Rightarrow *greater improvement*.
- Improvement also limited by the GET, i.e., the learned kernel is still “linear”.



$\sigma = \text{ReLU}, \sigma^* = \text{erf}$.

$\eta = \Theta(\sqrt{N})$ – Upper Bound via Nonparametric Analysis

Sufficiently large $\eta \Rightarrow \mathbf{W}_1$ travels far away from initialization.

☺ learned kernel can be “nonlinear”. ☹ Gaussian equivalence no longer holds.

Theorem ([BES+22] Upper Bound on Prediction Risk)

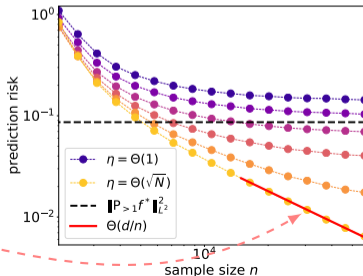
After one GD step on \mathcal{W} with $\eta = \Theta(\sqrt{N})$, for appropriate λ and $\psi_1 > \psi^*$,

$$\mathcal{R}_{\text{GD}}(\lambda) \leq 10\tau^* + \Theta(\psi_1^{-1}), \quad \text{w.h.p.},$$

where constant τ^* depends on σ, f^* , but not the specific value of step size η .

If $\tau^* \ll \|\mathbb{P}_{>1} f^*\|_{L^2}^2$, one feature learning step can outperform kernel lower bound:

- $\sigma = f^* = \tanh$: $\mathcal{R}_{\text{GD}}(\lambda) < \|\mathbb{P}_{>1} f^*\|_{L^2}^2$
- $\sigma = f^* = \text{erf}$: $\mathcal{R}_{\text{GD}}(\lambda) = O(d/n)$



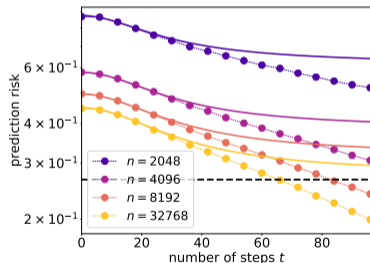
Conclusion and Future Directions

Random matrix theory allows us to characterize

- ❑ Precise conditions that determine the *sign of optimal ridge penalty*.
- ❑ Benefit of *representation learning* in the “early phase” of gradient descent.

Open Questions

- **Universality:** Under what conditions on the representation do we expect the RMT predictions to hold?
- **Beyond Universality:** What theoretical tools can we employ when the RMT predictions fail?



Failure case of RMT prediction.

Conclusion and Future Directions

Deep learning phenomena → interesting mathematical problems

- ❑ New models of (nonlinear) random matrix theory.
 - properties of neural net representation, beyond the proportional regime, ...
- ❑ What functions can be efficiently learned by neural network + gradient descent?
 - sparsity & low-dimensional structure, information exponent, ...
 - the role of architecture (depth, normalization, etc.) and optimization method (stochastic gradient, preconditioning, etc.)

Theoretical advances → principled guidance in practical settings

- ❑ How do we scale hyperparameters in the overparameterized setting?
 - selection of learning rate, regularization parameters, etc.
- ❑ “neural scaling laws” beyond the kernel regime.
 - How many samples, parameters, and optimization steps is required to achieve a desired test performance?

Conclusion

Thank you! Happy to take questions:)



Shun-ichi Amari



Jimmy Ba



Murat A. Erdogdu



Roger Grosse



Xuechen Li



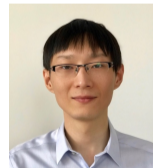
Atsushi Nitanda



Taiji Suzuki



Zhichao Wang



Ji Xu



Greg Yang

This Talk:

- [\[WX20\]](#), *Optimal weighted ℓ_2 regularization in overparameterized linear regression.*
- [\[ABG+21\]](#), *When does preconditioning help or hurt generalization?*
- [\[BES+22\]](#), *High-dimensional asymptotics of feature learning: how one gradient step improves the representation.*

Additional References:

- [\[NWS21\]](#), *Particle dual averaging: optimization of mean-field neural networks with global convergence rate analysis.*
- [\[NWS22\]](#), *Convex analysis of the mean-field Langevin dynamics.*
- [\[BES+23\]](#), *Learning in the presence of low-dimensional structure: a spiked random matrix perspective.*
- [\[MWS+23\]](#), *Gradient-based feature learning under structured data.*
- [\[SWN23\]](#), *Mean-field Langevin dynamics: time and space discretization, stochastic gradient, and variance reduction.*
- [\[SWO+23\]](#), *Feature learning via mean-field Langevin dynamics: classifying sparse parities and beyond.*