# High-dimensional Asymptotics of Feature Learning: How One Gradient Step Improves the Representation

Jimmy Ba[1], Murat A. Erdogdu[1], Taiji Suzuki[2], Zhichao Wang[3],
Denny Wu[1], Greg Yang[4]

[1]University of Toronto and Vector Institute
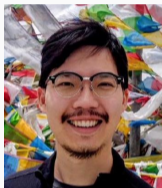[2]University of Tokyo and RIKEN AIP
[3]University of California, San Diego
[4]Microsoft Research AI

# Introduction

- <u>[BES+22]</u> Ba, Erdogdu, Suzuki, Wang, Wu, Yang. "*High-dimensional asymptotics of feature learning: how one gradient step improves the representation*".



Jimmy Ba



Murat A. Erdogdu



Taiji Suzuki



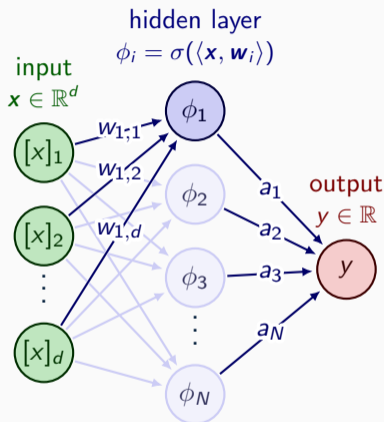Zhichao Wang



Denny Wu



Greg Yang

**Width-$N$ Two-layer NN**

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

- Input data: $\boldsymbol{x} \in \mathbb{R}^d$.
- Trainable parameters: $\boldsymbol{W} \in \mathbb{R}^{d \times N}, \boldsymbol{a} \in \mathbb{R}^N$.
- Element-wise nonlinearity: $\sigma : \mathbb{R} \to \mathbb{R}$.

**Optimization:** given a convex loss $\ell$,

- Optimizing $\boldsymbol{a}$ under fixed $\boldsymbol{W}$ is *convex*.
- Optimizing $\boldsymbol{W}$ under fixed $\boldsymbol{a}$ is *non-convex*.



hidden layer
$\phi_i = \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_i \rangle)$

input
$\boldsymbol{x} \in \mathbb{R}^d$

output
$y \in \mathbb{R}$

**Our Goal:** precise characterization of the performance of the trained NN.

3

- **Training.** Empirical risk minimization (potentially $\ell_2$-regularized):

$$\mathcal{L}(f) = \frac{1}{n} \sum_{i=1}^{n} (f(\mathbf{x}_i) - y_i)^2, \quad y_i = f^*(\mathbf{x}_i) + \varepsilon_i,$$

where $f^*$ is the target function (teacher model), and $\varepsilon$ is i.i.d. label noise.

- **Test.** Prediction risk: $\mathcal{R}(f) = \mathbb{E}_{\mathbf{x}}[(f(\mathbf{x}) - f^*(\mathbf{x}))^2] = \|f - f^*\|_{L^2(P_x)}^2$.

> **Regime of Interest** – **Proportional asymptotic limit**: $n, d, N \to \infty$,
> $n/d \to \psi_1$, $N/d \to \psi_2$, where $\psi_1, \psi_2 \in (0, \infty)$.
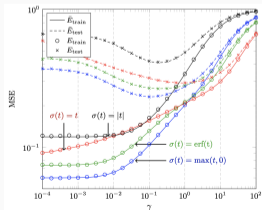
Why is this an interesting regime to analyze?

- It corresponds to the setting where the network width and data size are comparable, which is consistent with practical choices of model scaling.
- It might be possible to derive the *precise* prediction risk in this limit.
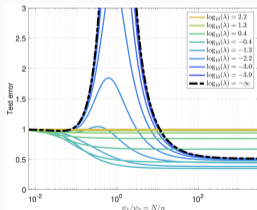
4

# Kernel Models Related to NN

Two widely-studied kernels derived from two-layer NN:

- **Conjugate Kernel (CK)** with features: $\phi_{\mathsf{CK}}(x) = \frac{1}{\sqrt{N}}\sigma(W^\top x) \in \mathbb{R}^N$.
  Regression on the CK corresponds to fixing $W$ and only learning the 2nd layer $a$.

- **Tangent Kernel (NTK)** with features: $\phi_{\mathsf{NT}}(x) = \frac{1}{\sqrt{Nd}}\mathsf{Vec}\big(\sigma'(W^\top x)x^\top\big) \in \mathbb{R}^{Nd}$.
  This kernel arises from gradient descent on certain wide neural networks.

When $W$ is randomly initialized, we arrive at a **random features (RF)** model, the precise asymptotics of which has been extensively studied in the proportional limit.



[Louart, Liao, and Couillet, 2018].



[Mei and Montanari, 2019].

## Limitation of Kernel Ridge Regression

Can these RF models fully capture the effectiveness of NNs? *Not quite...*

Consider the *ridge regression estimator* for $RF \in \{CK, NT\}$:

$$f_{RF}^\lambda(\boldsymbol{x}) = \langle \phi_{RF}(\boldsymbol{x}), \hat{\boldsymbol{a}}_\lambda \rangle, \quad \hat{\boldsymbol{a}}_\lambda = \text{argmin}_{\boldsymbol{a}} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \langle \phi_{RF}(\boldsymbol{x}_i), \boldsymbol{a} \rangle)^2 + \frac{\lambda}{N} \|\boldsymbol{a}\|_2^2 \right\}.$$

---

**Theorem (Ghorbani et al. 19, Hu and Lu 20, Bartlett et al. 21, ...)**

[Informal] *Denote $P_{>1}$ as the projector orthogonal to constants and linear functions in $L^2(P_X)$. Then under certain concentration conditions on the input $\boldsymbol{x}$, we have[1]*

$$\inf_{\lambda > 0} \min \left\{ \mathcal{R}(f_{CK}^\lambda), \mathcal{R}(f_{NT}^\lambda) \right\} \geq \boxed{\|P_{>1} f^*\|_{L^2}^2} + o_{d, \mathbb{P}}(1),$$

---

- In the proportional limit, RF models can only learn **linear functions**.
- NNs are clearly more powerful than linear models on the input...

---

[1] Similar lower bound also holds for certain rotationally invariant kernels studied in [El Karoui 10].

# Feature Learning in Two-layer NN

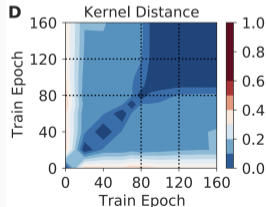Where does this gap come from?   **Feature Learning!**

- When we optimize the first-layer parameters $W$, we expect the model to "adapt" to the data and learn useful representations.
- In RF models, $W$ is fixed, so there is no "representation learning".

**Motivation:** Can we precisely capture the presence of *feature learning* in the proportional limit, when the first-layer $W$ is optimized via *gradient descent*?

**Empirical Observation:**
- Neural network features often change most rapidly in the **early phase** of gradient descent (GD) training.

We consider the most simplified setting of the "early phase": one gradient step on $W$, and analyze how the learned CK adapts to the learning problem.
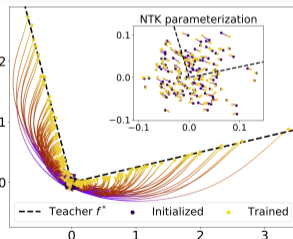


[Fort et al. 2020].

# Problem Setting: Basic Assumptions

1. **Proportional Limit.** $n, d, N \to \infty$, $n/d \to \psi_1$, $N/d \to \psi_2$, $\psi_1, \psi_2 \in (0, \infty)$.

2. **Student-teacher Setup.** $y_i = f^*(x_i) + \varepsilon_i$, where $x_i \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, I)$, $\varepsilon_i$ is i.i.d. noise with variance $\sigma_\varepsilon^2$, and $f^*$ is Lipschitz with $\|f^*\|_{L^2} = \Theta_d(1)$.

3. **Normalized Activation.** $\sigma$ has bounded first three derivatives, and is normalized such that $\mathbb{E}[\sigma(z)] = 0$, $\mathbb{E}[z\sigma(z)] = \mu_1 \neq 0$, for $z \sim \mathcal{N}(0, 1)$.

4. **Gaussian Initialization.** $[W_0]_{ij} \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1/d)$, $[a]_j \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1/N)$.

**Note:** we use the <u>mean-field</u> parameterization[2], which admits a *feature learning limit* (i.e., the weights do not "freeze" around the initialization).

$$f(x) = \frac{1}{\sqrt{N}} \sum_{i=1}^{N} a_i \sigma(\langle x, w_i \rangle) = \underbrace{\frac{1}{\sqrt{N}} a^\top}_{\approx 1/N} \sigma(W^\top x).$$



NNs trained till $\mathcal{L}(f) < 10^{-3}$.

---

[2] The NTK scaling corresponds to dropping the $\frac{1}{\sqrt{N}}$-prefactor.

## Problem Setting: One-step Gradient Descent

- **One-step GD on 1st Layer**. We take **one gradient step**[3] on the empirical MSE loss $\mathcal{L}(f) = \frac{1}{n}\sum_{i=1}^{n}(f(\boldsymbol{x}_i) - y_i)^2$, that is, $\boldsymbol{W}_1 = \boldsymbol{W}_0 + \eta\sqrt{N}\cdot\boldsymbol{G}_0$, where

$$\boldsymbol{G}_0 := -\frac{1}{n}\boldsymbol{X}^{\top}\left[\left(\frac{1}{\sqrt{N}}\left(\frac{1}{\sqrt{N}}\sigma(\boldsymbol{X}\boldsymbol{W}_0)\boldsymbol{a} - \boldsymbol{y}\right)\boldsymbol{a}^{\top}\right)\odot\sigma'(\boldsymbol{X}\boldsymbol{W}_0)\right],$$

- **Ridge Regression for 2nd Layer**. After learning the features for one step, we perform ridge regression on the trained CK using **a fresh set of data** $\{\tilde{\boldsymbol{X}},\tilde{\boldsymbol{y}}\}$:

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

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

**This Work**: We aim to compute $\mathcal{R}_{\mathrm{GD}}(\lambda)$, and show its *improvement* over the initialized RF, and potentially over the *lower bound* $\|\mathrm{P}_{>1}f^*\|_{L^2}^2$.

**Challenge:** cannot directly use *random matrix theory*, as $\boldsymbol{W}_1$ is no longer "random".

---

[3]Some of our results also apply to multiple gradient steps on $\boldsymbol{W}$.

# Properties of the Gradient Matrix $G_0$

Can we exploit certain structure of the first GD step to simplify the calculation?

**Orthogonal Decomposition of $\sigma$:**
$\sigma(z) = \mu_1 z + \sigma_\perp(z)$, where $\mu_1 = \mathbb{E}[\sigma'(z)] \quad \Rightarrow \quad \mathbb{E}[\sigma_\perp(z)] = \mathbb{E}[z\sigma_\perp(z)] = 0$.

---

**Proposition (BES+22)**

Recall $G_0 = \frac{1}{\eta\sqrt{N}}(W_1 - W_0)$. Define rank-1 matrix $A := \frac{\mu_1}{n\sqrt{N}} X^\top y a^\top$. Then

$$\sqrt{N} \cdot \|G_0 - A\| \lesssim \|G_0\|, \text{ w.h.p.}$$

---

**Intuition:** Many commonly-used activations are monotone, so $\sigma'$ is not centered:

$$n\sqrt{N} \cdot G_0 = \mu_1 X^\top (y - f_0(X)) a^\top + X^\top ((y - f_0(X)) a^\top \odot \sigma'_\perp(X W_0))$$

Hence $G_0$ contains:  $\qquad \|A\|_F \asymp \|B\|_F$, but $\|A\| \gg \|B\|$

- A rank-1 "spike" $A$
- A "residual" with smaller operator norm (but not Frobenius norm) $B$

## Selection of Learning Rate $\eta$

Based on the decomposition of $\boldsymbol{G}_0$, we focus on the following choices[4] of $\eta$:

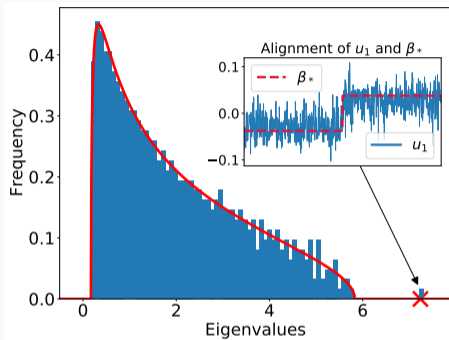> - <u>Small lr</u>: $\eta = \Theta(1) \Rightarrow \|\boldsymbol{W}_1 - \boldsymbol{W}_0\| \asymp \|\boldsymbol{W}_0\|$.
> - <u>Large lr</u>: $\eta = \Theta(\sqrt{N}) \Rightarrow \|\boldsymbol{W}_1 - \boldsymbol{W}_0\|_F \asymp \|\boldsymbol{W}_0\|_F$.

**Remarks:**

- Under $\eta = \Theta(1)$, the NN after one GD step remains close to the **kernel regime**: each neuron (or parameter) does not travel far away from the initialization, i.e., $\left|[\boldsymbol{W}_1 - \boldsymbol{W}_0]_{ij}\right| \ll \left|[\boldsymbol{W}_0]_{ij}\right|$ for all $i, j$ with high probability.

- $\eta = \Theta(\sqrt{N})$ mirrors the **maximal update parameterization** [Yang and Hu 2020]: for $\boldsymbol{x} \sim \mathcal{N}(0, \boldsymbol{I})$, the change in each coordinate of the feature vector is significant, i.e., $\left|\sigma(\boldsymbol{W}_1^\top \boldsymbol{x}) - \sigma(\boldsymbol{W}_0^\top \boldsymbol{x})\right|_i \asymp \left|\sigma(\boldsymbol{W}_0^\top \boldsymbol{x})\right|_i = \tilde{\Theta}(1)$ for all $i$ with high probability.

---

[4]For smaller $\eta = o(1)$, one can easily verify that change in the prediction risk is negligible.

# A Spiked Model for $W_1$



Alignment of $u_1$ and $\beta_*$

**Blue**: empirical simulation.
**Red**: analytic prediction.
*(next slide)*

- $\sigma = \tanh$, $f^*(x) = \text{ReLU}(\langle x, \beta_* \rangle)$.
- Teacher vector $\beta_* \propto [-1_{d/2}; 1_{d/2}]$.
- $\psi_1 = n/d = 4$, $\psi_2 = N/d = 2$.
- $\eta = 2$.

**Observation:** after one gradient step with learning rate $\eta = \Theta(1)$:

- The **bulk** of the spectrum of $W$ remains unchanged [5].

- A **spike** ($\times$) appears in $W_1$, which aligns with linear component of $f^*$.

---

[5] The spectrum of the initialized $W_0$ is characterized by the Marchenko–Pastur law.

**Orthogonal Decomposition:** $f^*(\boldsymbol{x}) = \mu_0^* + \mu_1^* \langle \boldsymbol{x}, \boldsymbol{\beta}_* \rangle + \mathsf{P}_{>1} f^*(\boldsymbol{x})$,

• *Linear part:* $\|\boldsymbol{\beta}_*\| = 1$, $\mu_1^* \boldsymbol{\beta}_* = \mathbb{E}[\boldsymbol{x} f^*(\boldsymbol{x})]$; • *Noninear part:* $\|\mathsf{P}_{>1} f^*\|_{L^2} = \mu_2^*$.
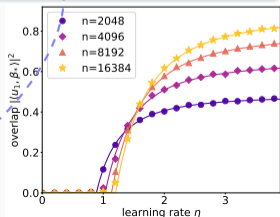
---

**Theorem (BES+22)**

*For $\eta = \Theta(1)$, define $\theta_1 := \sqrt{\|f^*\|_{L^2}^2 \psi_1^{-1} + \mu_1^{*2}} \cdot \mu_1 \eta$, $\theta_2 := \mu_1 \mu_1^* \eta$. The leading singular value $s_1(\boldsymbol{W}_1)$ and the corresponding singular vector $\boldsymbol{u}_1$ satisfy*

$$s_1(\boldsymbol{W}_1) \to \sqrt{\frac{(1+\theta_1^2)(\psi_2 + \theta_1^2)}{\theta_1^2}}, \qquad |\langle \boldsymbol{u}_1, \boldsymbol{\beta}_* \rangle|^2 \to \frac{\theta_2^2}{\theta_1^2} \left( 1 - \frac{\psi_2 + \theta_1^2}{\theta_1^2(\theta_1^2 + 1)} \right),$$

*for $\theta_1 > \psi_2^{1/4}$; otherwise, $s_1(\boldsymbol{W}_1) \to 1 + \sqrt{\psi_2}$, $|\langle \boldsymbol{u}_1, \boldsymbol{\beta}_* \rangle| \to 0$.*

When $\eta$ exceeds some threshold, a "spike" appears:

• Increase step size $\eta$ $\Rightarrow$ <u>larger spike</u> $s_1(\boldsymbol{W}_1)$.

• Increase sample size $\psi_1$ $\Rightarrow$ <u>greater alignment</u>.

## A Spiked Model for CK?

**Question:** How does the spike in $W_1$ affect the *kernel (CK) matrix*?

For $\eta = \Theta(1)$, and *odd activation* $\sigma$, the <u>expected</u> CK matrix $\Sigma_\Phi$ satisfies

$$\left\| \Sigma_\Phi - \overline{\Sigma}_\Phi \right\| \xrightarrow{\mathbb{P}} 0, \text{ where } \Sigma_\Phi = \mathbb{E}_x \left[ \sigma(W_1^\top x) \sigma(x^\top W_1) \right], \ \overline{\Sigma}_\Phi = \mu_1^2 W_1^\top W_1 + \mu_2^2 I.$$

- Intuitively, we expect a spike to appear in the (empirical) CK matrix.

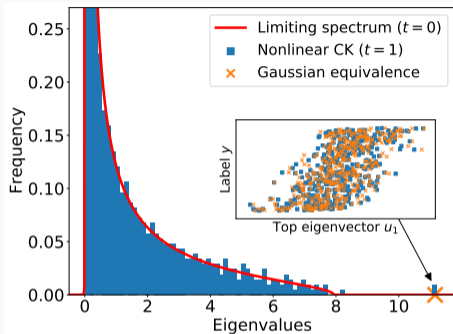- How do we predict properties of the CK spike?   **Gaussian Equivalence**

$\boxed{\text{Nonlinear CK}} : \Phi = \frac{1}{\sqrt{N}} \sigma(\tilde{X} W_1), \quad \boxed{\text{"Linearized" CK}} : \bar{\Phi} = \frac{1}{\sqrt{N}} \left( \mu_1 \tilde{X} W_1 + \mu_2 Z \right).$

**Conjecture (Gaussian Equivalence of CK Spike)**

*For odd activation $\sigma$ and $\eta = \Theta(1)$, given i.i.d. training data $\tilde{X}, \tilde{y}$ (independent to $W_1$). Denote the left singular vectors of $\Phi, \bar{\Phi}$ as $u_1, \bar{u}_1$, we conjecture*

$$\left| s_i(\Phi) - s_i(\bar{\Phi}) \right| = o_{d,\mathbb{P}}(1), \ \forall i \in [n]; \quad |\langle u_1, \tilde{y}/\|\tilde{y}\| \rangle|^2 = |\langle \bar{u}_1, \tilde{y}/\|\tilde{y}\| \rangle|^2 + o_{d,\mathbb{P}}(1).$$

**Blue**: empirical simulation.
**Red**: analytic prediction (initial CK).
**Orange**: Gaussian equivalence.

- $\sigma = $ SoftPlus.
- $f^*(x) = \tanh(\langle x, \beta_* \rangle)$.
- $\psi_1 = n/d = 3/2$, $\psi_2 = N/d = 5/4$.
- $\eta = 2$.

- The **bulk** of the CK spectrum remains unchanged [6].

- A **spike** ($\times$) appears in the learned CK, predicted by Gaussian equivalence.

- The corresponding eigenvector $u_1$ aligns with training labels $\tilde{y}$.

---

[6] The spectrum of the initialized $CK_0$ is characterized in [Fan and Wang 2020].

# Prediction Risk of CK Ridge Regression

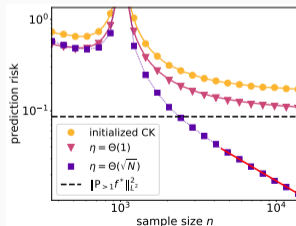**Question:** does this alignment improve the performance of the kernel model?

> **Case Study: Single-index target[7].** $f^*(x) = \sigma^*(\langle x, \beta^* \rangle)$.
> where $\|\beta_*\| = 1$, and $\sigma^*$ is Lipschitz with $\mu_0^* = 0$, $\mu_1^* \neq 0$.

**Goal:** compute the prediction risk $\mathcal{R}_{\mathrm{GD}}(\lambda)$ of the ridge estimator

$$f_{\mathrm{GD}}^\lambda(x) = \frac{1}{\sqrt{N}} \hat{a}_\lambda^\top \sigma\left(W_1^\top x\right), \ \hat{a}_\lambda = \mathrm{argmin}_a \left\{ \frac{1}{n} \left\| \tilde{y} - \frac{1}{\sqrt{N}} \sigma(\tilde{X} W_1) a \right\|^2 + \frac{\lambda}{N} \|a\|^2 \right\}.$$

We consider the following learning rate scalings:



- Small lr $\eta = \Theta(1)$ : trained CK <u>always improve</u> upon the initial CK ridge estimator ($\mathcal{R}_0(\lambda)$).

- Large lr $\eta = \Theta(\sqrt{d})$ : for some $f^*$, trained CK may outperform the <u>lower bound</u> $\|\mathsf{P}_{>1} f^*\|_{L^2}$.

---

[7]This setting is often studied in RF regression (e.g. [Gerace et al. 20],[Dhifallah and Lu 20]).

# The Gaussian Equivalence Property

Consider the prediction risk of ridge regression on features $F \in \{CK, GE\}$:

$$\mathcal{R}_F(\lambda) = \mathbb{E}_x \big(\langle \phi_F(x), \hat{a}_\lambda \rangle - f^*(x)\big)^2, \ \hat{a}_\lambda = \operatorname{argmin}_a \Big\{ \frac{1}{n} \sum_{i=1}^n (y_i - \langle \phi_F(x_i), a \rangle)^2 + \frac{\lambda}{N} \|a\|^2 \Big\}$$

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

  where $\mu_1 = \mathbb{E}[z\sigma(z)]$, $\mu_2 = \sqrt{\mathbb{E}[\sigma(z)^2] - \mu_1^2}$

The **Gaussian Equivalence Property** refers to: $\mathcal{R}_{CK}(\lambda) \approx \mathcal{R}_{GE}(\lambda)$.

Previously, the Gaussian equivalence theorem (GET) has been shown for certain <u>RF models</u> [Hu and Lu 2020], but not for the <u>trained features</u>.

**Implications of the Gaussian Equivalence:**

- We can equivalently compute $\mathcal{R}_{GE}$, which can be handled via RMT tools ☺
- The <u>nonlinear</u> CK model achieves the same performance as a <u>linear</u> model ☹

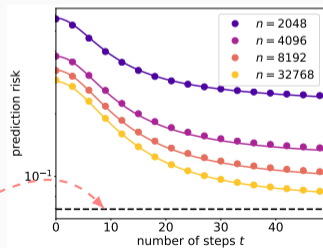# Gaussian Equivalence for Trained Features

**Theorem (BES+22)**

*Assume $\sigma$ is **odd** in addition to the previous assumptions, then for fixed $t \in \mathbb{N}$, after the first-layer $\boldsymbol{W}$ is trained for t gradient steps with $\eta = \Theta(1)$,*
$$|\mathcal{R}_{\mathrm{CK}}(\lambda) - \mathcal{R}_{\mathrm{GE}}(\lambda)| = o_{d,\mathbb{P}}(1), \text{ for } \lambda > 0.$$

**Intuition:** GET holds when $\boldsymbol{W}_t$ is not far away from the random initialization $\boldsymbol{W}_0$.

**Figure**: dots represent empirical values, solid curves are asymptotics predicted by CGMT.

- For learning rate $\eta = \Theta(1)$, GET remains accurate in the early phase of training

- Prediction risk $\mathcal{R}_{\mathrm{GD}}(\lambda)$ can improve, but is still lower-bounded by $\|\mathrm{P}_{>1}f^*\|_{L^2}^2$



$\sigma = \mathsf{ReLU}, \sigma^* = \mathsf{tanh}.$

## Gaussian Equivalence Theorem (continued)

**Proof Sketch.** We extend the argument in [Hu and Lu 2020] outline below.

1. **Lindeberg exchange**. Let $\hat{\boldsymbol{g}}_k$ be the solution of the optimization problem:

$$L_k \triangleq \min_{\boldsymbol{g} \in \mathbb{R}^N} \left\{ \sum_{i=1}^{k} \ell(y_i, \langle \boldsymbol{g}, \phi_{\mathrm{GE}}(\boldsymbol{x}_i)\rangle) + \sum_{j=k+1}^{n} \ell(y_j, \langle \boldsymbol{g}, \phi_{\mathrm{CK}}(\boldsymbol{x}_j)\rangle) + \frac{n}{N}\left(\lambda\|\boldsymbol{g}\|_2^2 + Q(\boldsymbol{g})\right) \right\}$$

As there are $N$ *total swaps*, it suffices to show that for bounded test function $\zeta$,

$$\left|\mathbb{E}\zeta\left(\tfrac{1}{N}L_k\right) - \mathbb{E}\zeta\left(\tfrac{1}{N}L_{k-1}\right)\right| = \mathcal{O}\left(\frac{\mathrm{polylog}N}{N^{3/2}}\right). \tag{A}$$

2. **Central limit theorem**. A crucial step in establishing (A) is the following CLT:

$$\left|\mathbb{E}\varphi(\langle\phi_{\mathrm{GE}}, \boldsymbol{g}\rangle) - \mathbb{E}\varphi(\langle\phi_{\mathrm{CK}}, \boldsymbol{g}\rangle)\right| = \mathcal{O}\left(\frac{\mathrm{polylog}N}{\sqrt{N}} \cdot \left(1 + \|\boldsymbol{g}\|_\infty^2\right)\right).$$

This is shown using *Stein's method*, when $\boldsymbol{W}$ has *near-orthogonal* columns.

3. $\ell_\infty$-**norm control**. Finally, we show that entries of $\hat{\boldsymbol{g}}_k$ are "underline{evenly distributed}"[8]:

$$\mathbb{P}\left(\|\hat{\boldsymbol{g}}_k\|_\infty \geq \mathrm{polylog}N\right) \leq \exp\left(-c\log^2 N\right), \text{ for all } k \in [N].$$

---

[8] In this part of the analysis, [Hu and Lu 2020] required $\boldsymbol{W}_{ij}$ to be i.i.d. Gaussian.

**Goal:** can we rigorously show that one feature learning step always *decreases* the prediction risk of the CK ridge regression estimator?

- Risk of *initial* CK (random features): $\mathcal{R}_0(\lambda) = \mathbb{E}_{\boldsymbol{x}}\big(\langle\sigma(\boldsymbol{W}_0^\top \boldsymbol{x}),\hat{\boldsymbol{a}}_0\rangle - f^*(\boldsymbol{x})\big)^2$.

- Risk of *trained* CK (after one step): $\mathcal{R}_{\mathrm{GD}}(\lambda) = \mathbb{E}_{\boldsymbol{x}}\big(\langle\sigma(\boldsymbol{W}_1^\top \boldsymbol{x}),\hat{\boldsymbol{a}}_1\rangle - f^*(\boldsymbol{x})\big)^2$.

### Theorem (BES+22)

*For $\eta = \Theta(1)$ and $\lambda > 0$, as $n/d \to \psi_1$, $N/d \to \psi_2$, we have*

$$\mathcal{R}_0(\lambda) - \mathcal{R}_{\mathrm{GD}}(\lambda) \xrightarrow{\mathbb{P}} \delta(\eta, \lambda, \psi_1, \psi_2).$$

- $\delta(\eta, \lambda, \psi_1, \psi_2)$ *is a **non-negative** function of $\eta, \lambda, \psi_1, \psi_2 \in (0, +\infty)$;*
- $\delta$ *vanishes if and only if (at least) one of $\mu_1^*, \mu_1$ and $\eta$ is zero.*
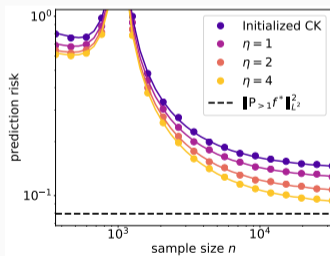
Provable improvement over the initial CK model!

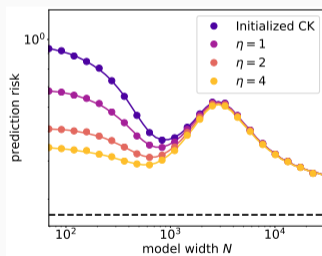**Note:** this does not require the student and teacher to have the same nonlinearity

In some special cases, the expression of $\delta$ can be further simplified.

**Proposition (BES+22)**

- **[Large sample limit]** As $\psi_1 \to \infty$, $\delta$ is ***increasing*** with respect to $\eta$.
- **[Large width limit]** As $\psi_2 \to \infty$, $\delta(\eta, \lambda, \psi_1, \psi_2) \to 0$.



Risk vs. sample size.



Risk vs. model width.

**Note:** In all cases, $\mathcal{R}_0(\lambda) \geq \mathcal{R}_{\mathrm{GD}}(\lambda) \geq \|P_{>1}f^*\|_{L^2}^2$ due to the GET under $\eta = \Theta(1)$.

Finally, we consider the <u>large learning rate</u> regime with $\eta = \Theta(\sqrt{d})$.

- $W_1$ travels far away from initialization $\Rightarrow$ CK can be "nonlinear" ☺
- In the absence of GET, precise analysis of prediction risk is difficult ☹

**Alternative:** *upper-bound* $\mathcal{R}_{\mathrm{GD}}(\lambda)$ and compare against *kernel lower bound*.

$$\text{We define: } \tau^* := \inf_\eta \mathbb{E}_{\xi_1} \left( \sigma^*(\xi_1) - \mathbb{E}_{\xi_2}(\sigma(\eta\xi_1 + \xi_2)) \right)^2$$

---

**Lemma (BES+22)**

[Informal] *Given **bounded** activation $\sigma$, after one GD step on $W$ with $\eta = \Theta(\sqrt{N})$, there exists some $\tilde{f}(x) = \frac{1}{\sqrt{N}} \tilde{a}^\top \sigma(W_1^\top x)$ that achieves prediction risk "close" to $\tau^*$.*

---

- $\tau^*$ can be interpreted as some measure of "model misspecification".
- **Note:** the definition of $\tau^*$ does not involve the specific value of step size $\eta$.
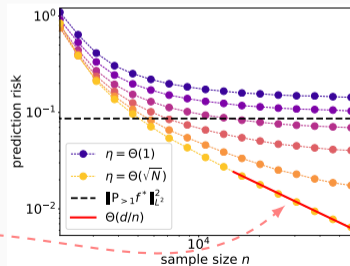
**Theorem (BES+22)**

*After one GD step on $W$ with $\eta = \Theta(\sqrt{N})$, there exist constants $C, \psi_1^* > 0$ such that for any $\psi_1 > \psi_1^*$, and $n^{\epsilon-1} < N^{-1}\lambda < n^{-\epsilon}$ for some small $\epsilon > 0$, we have*
$$\mathcal{R}_{\mathrm{GD}}(\lambda) \leq 16\tau^* + C\left(\sqrt{\tau^*} \cdot \psi_1^{-1/2} + \psi_1^{-1}\right),$$
*with probability 1, as $n, d, N \to \infty$ proportionally.*

If $\tau^* \ll \|P_{>1}f^*\|_{L^2}^2$, CK ridge regression after *one feature learning step* outperforms the kernel ridge lower bound:

- $\underline{\sigma = \sigma^* = \tanh}$: $\mathcal{R}_{\mathrm{GD}}(\lambda) < \|P_{>1}f^*\|_{L^2}^2$

- $\underline{\sigma = \sigma^* = \mathrm{erf}}$: there exists constant $C > 0$ s.t. $\mathcal{R}_{\mathrm{GD}}(\lambda) \leq C \cdot \psi_1^{-1} = \Theta(d/n)$



**Caution:** separation only present in specific $(\sigma, \sigma^*)$

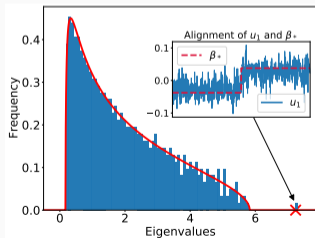$\sigma = \sigma^* = \mathrm{erf}$, $\eta = N^\alpha$, $\alpha \in [0, 1/2]$.

**How Does One Gradient Step Change the Weights?**

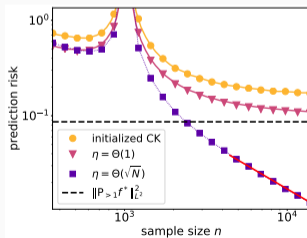- The isolated singular vector of $W_1$ aligns with *linear component of $f^*$*.
- The top eigenvector of CK matrix aligns with *training labels $y$* (conjecture).

**How Do the Learned Features Improve Generalization?**

- $\eta = \Theta(1)$ – **Linear Regime**. Precise analysis via GET; $\mathcal{R}_0 \geq \boxed{\mathcal{R}_{\mathrm{GD}}} \geq \|\mathrm{P}_{>1}f^*\|_{L^2}^2$.
- $\eta = \Theta(\sqrt{d})$ – **Nonlinear Regime**. For certain $f^*$, $\mathcal{R}_0 \geq \|\mathrm{P}_{>1}f^*\|_{L^2}^2 \geq \boxed{\mathcal{R}_{\mathrm{GD}}}$.
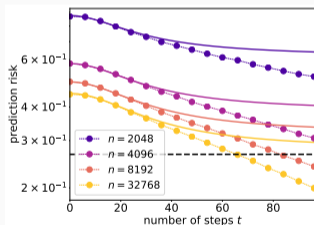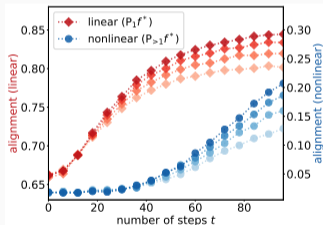


Signal+noise structure of $W_1$.



Improvement of prediction risk.

**Some questions to consider:**

1. A *spiked model* for the kernel (CK) matrix after one gradient step?

2. "Phase transition" in the Gaussian equivalence property?

3. *Precise asymptotics* beyond Gaussian equivalence?



Prediction risk vs. time step $t$.



Alignment with teacher $f^*$.

## Thank you!

# References

- El Karoui, 2010. *The spectrum of kernel random matrices*.

- Louart, Liao, and Couillet, 2018. *A random matrix approach to neural networks*.

- Mei and Montanari, 2019. *The generalization error of random features regression: Precise asymptotics and double descent curve*.

- Ghorbani et al., 2019. *Linearized two-layer neural networks in high dimensions*.

- Hu and Lu, 2020. *Universality laws for high-dimensional learning with random features*.

- Fan and Wang, 2020. *Spectra of the conjugate kernel and neural tangent kernel for linear-width neural networks*.

- Gerace et al., 2020. *Generalisation error in learning with random features and the hidden manifold model*.

- Yang and Hu, 2021. *Feature learning in infinite-width neural networks*.

- Loureiro et al., 2021. *Learning curves of generic features maps for realistic datasets with a teacher-student model*.

- Bartlett et al., 2021. *Deep learning: a statistical viewpoint*.