# Learning Theory of Transformers: Generalization and Optimization of In-Context Learning

## Taiji Suzuki

The University of Tokyo / AIP-RIKEN
(Deep learning theory team)

THE UNIVERSITY OF TOKYO

AIP

4th/Dec/2024

Simons Institute Workshop

## LLM (e.g., GPT3/GPT4)



[Alammar: How GPT3 Works - Visualizations and Animations, https://jalammar.github.io/how-gpt3-works-visualizations-animations/]

[Brown et al. "Language Models are Few-Shot Learners", NeurIPS2020]

## • Transformer



Figure 1: The Transformer - model architecture.
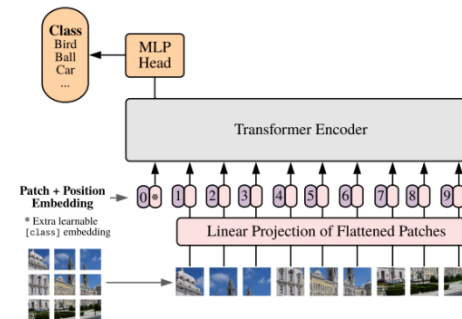
[Vaswani et al.: Attention is All you Need. NIPS2017]

## Vision tasks (e.g., ViT)



[Dosovitskiy et al.: An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. arXiv:2010.11929. ICLR2021]

Please explain what is ChatGPT

ChatGPT is an artificial intelligence (AI) language model created by OpenAI. It is designed to understand and generate human-like text based on the input it receives. Here's a detailed breakdown of what ChatGPT is and what it does:

## 1. Purpose

- ChatGPT is primarily used for conversation and can respond to questions, provide explanations, write creatively, assist with tasks, and more.

- It serves as a tool for various applications, such as customer support, learning assistance, idea brainstorming, and general information retrieval.

## 2. How It Works

- **Training:** ChatGPT is trained using a machine learning technique called reinforcement learning from human feedback (RLHF), which refines its ability to respond in helpful, accurate, and safe

[OpenAI, ChatGPT 4o]

# Why does it work so well?

Pretrained Large Language Models (LLMs) have significant ability of In-Context Learning (ICL) [Brown et al., 2020].

Please guess the number that fits in the '?'.

**context**

1,1 -> 2
2,3 -> 5
8,13 -> 21
6,0 -> 6
10,1 -> 11
5,27 -> ?

The pattern in the given pairs of numbers appears to be the sum of the two numbers.

So, the number that fits in the '?' is 32.

**Question**

**ChatGPT**

Pretrained Large Language Models (LLMs) have significant ability of In-Context Learning (ICL) [Brown et al., 2020].

Please guess the word that fits in the '?'.

left -> right
dark -> light
short -> long
small -> big
up -> ?

context

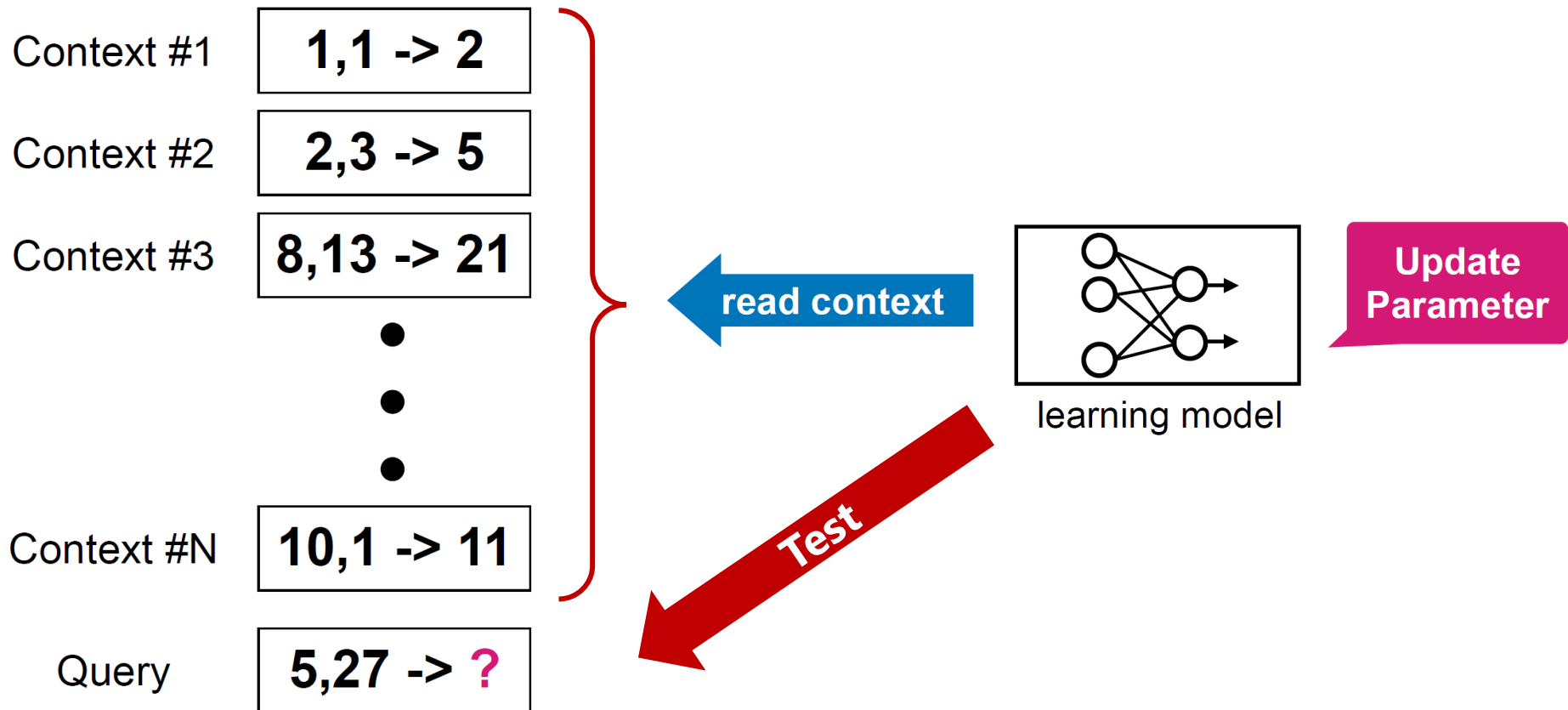The pattern in the given pairs of words seems to be antonyms:

So, the word that fits in the '?' is "down".

**Question**

**ChatGPT**

## Traditional "fine tuning" approach
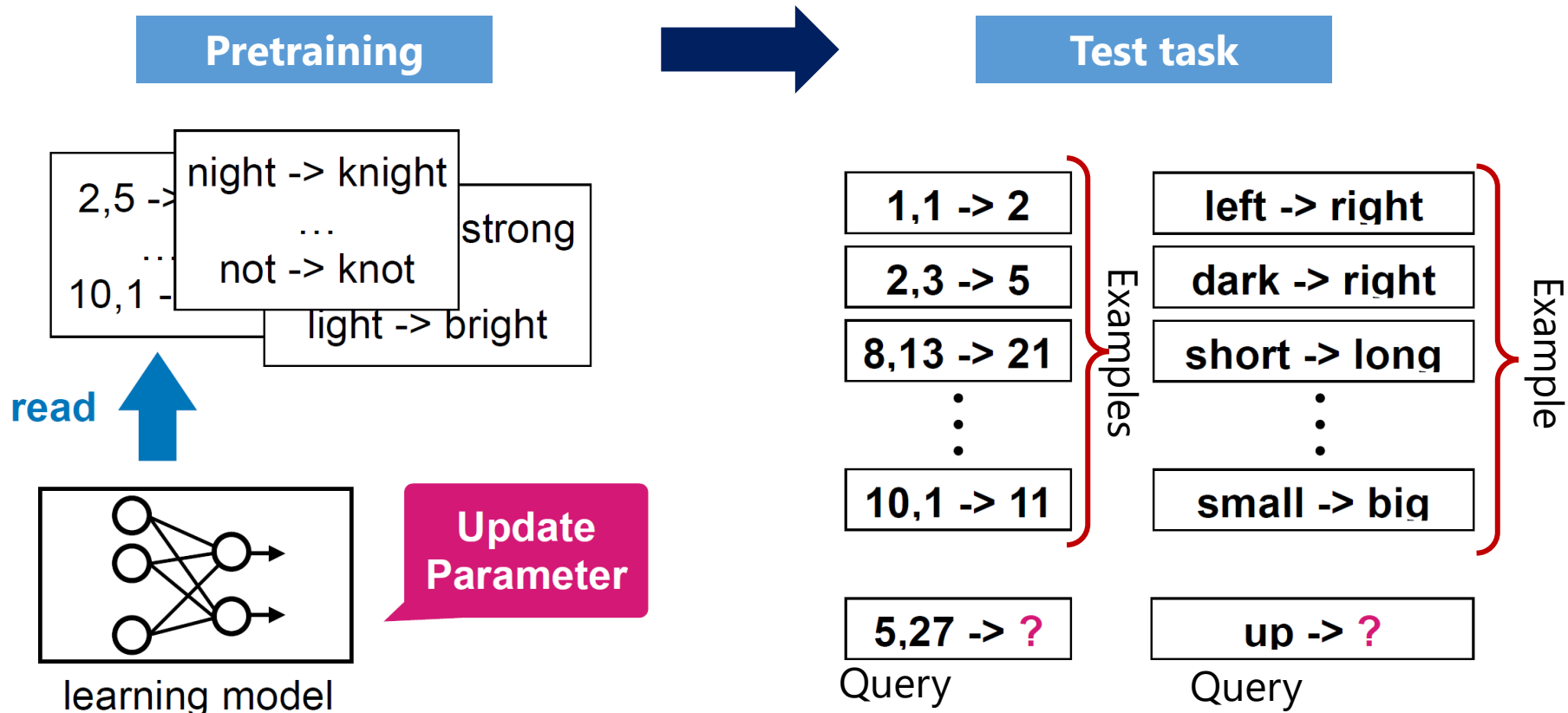


(e.g., RLHF)

ICL is performed **without updating model parameters** unlike the traditional "fine-tuning" regime in the test task.
→ Meta-learning



| Pretraining | Test task |

night -> knight
…
not -> knot

2,5 ->
…
10,1 -

strong
light -> bright

**read**

**Update Parameter**

learning model

1,1 -> 2
2,3 -> 5
8,13 -> 21
⋮
10,1 -> 11

Examples

left -> right
dark -> right
short -> long
⋮
small -> big

Example

5,27 -> **?**
Query

up -> **?**
Query

During pretraining, several tasks are observed to train the model.
→ Task generalization.

**Question:**
What mechanism allows a Transformer to perform ICL?

# Presentation overview

**Statistics**

**Minimax optimality**
- Nonparametric analysis
- Approximation error analysis

**Optimization**

**Global optimality of nonlinear feature learning**
- Mean field limit
- Strict saddle

**Statistics/Optimization**

**Feature learning with one step GD**
- Single index model
- Information exponent
- Advantage of pre-training

- [Minimax optimality and approximation error bound] Kim, Nakamaki, Suzuki: Transformers are Minimax Optimal Nonparametric In-Context Learners. NeurIPS2024
- [Optimization in mean field limit] Kim, Suzuki: Transformers Learn Nonlinear Features In Context: Nonconvex Mean-field Dynamics on the Attention Landscape. ICML2024 (arXiv:2402.01258).
- [Identifying low dimensional subspace with information exponent k] Oko, Song, Suzuki, Wu: Transformer efficiently learns low-dimensional functions in context. NeurIPS2024.

# Approximation theory/ Statistical analysis

# Nonparametric analysis of in-context learning

[Kim, Nakamaki, Suzuki: Transformers are Minimax Optimal Nonparametric In-Context Learners. NeurIPS2024]



**Juno Kim**

# Mathematical formulation of in-context learning

**Model:** $y_{i,t} = F_t^\circ(x_{i,t}) + \epsilon_{i,t} \qquad (i = 1, \ldots, n)$

$t = 1, \ldots, T$: Task index

- The true functions $F_t^\circ$ are different across different tasks.
- $F_t^\circ$ is generated randomly for each task.

**Pretraining ($T$ tasks) :**

$$X_t = [x_{1,t}; \ldots; x_{n,t}] \quad x_{\mathrm{qr},t}$$
$$\ldots \qquad \times T$$
$$Y_t = [y_{1,t}; \ldots; y_{n,t}] \quad y_{\mathrm{qr},t}$$

➤ We observe pretraining task data $T$ times.
➤ Each task has $n$ data.

**Test task (In-context learning) :**

$$X_{T+1} = [x_{1,T+1}; \ldots; x_{n,T+}$$
$$\ldots$$
$$Y_{T+1} = [y_{1,T+1}; \ldots; y_{n,T+1}]$$

**(Implicit) Bayes estimation**
➤ Learn prior at pretraining
➤ Perform posterior inference at the test task

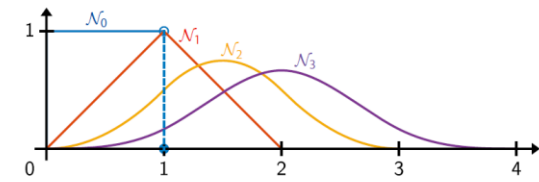Suppose that the true function admits a basis function decomposition:

$$F_t^\circ(x) = \beta_t^\top f^\circ(x)$$

where $\beta_t \sim (0, \Sigma)$ and $f^\circ(x) \in \mathbb{R}^\infty$.

---

- B-Spline (Besov)

$$\mathcal{N}(x) = \begin{cases} 1 & (x \in [0,1]), \\ 0 & (\text{otherwise}) \end{cases}$$

$$\mathcal{N}_m(x) = (\underbrace{\mathcal{N} * \mathcal{N} * \cdots * \mathcal{N}}_{m+1 \text{ times}})(x)$$



Tensor product B-spline:

$$M_{a,b}^d(x) = \prod_{j=1}^d \mathcal{N}_m(2^{a_j} - b_j)$$

$$f_j^\circ(x) = M_{a(j),b(j)}^d(x) \implies F_\beta^\circ = \beta^\top f^\circ \in B_{2,2}^\alpha$$

- Fourier (Sobolev, $\gamma$-smooth)

$$f_j^\circ(x) = \prod_{k=1}^\infty \sqrt{2}\cos(2\pi 2^{s_{j,k}} x_k - \delta_{j,k}\pi/2) \implies F_\beta^\circ \in \mathcal{F}_{2,2}^\gamma([0,1]^\infty)$$

$\gamma$-smooth function class for $d = \infty$ [Okumoto&Suzuki,ICLR2022], [Takakura&Suzuki, ICML2024]
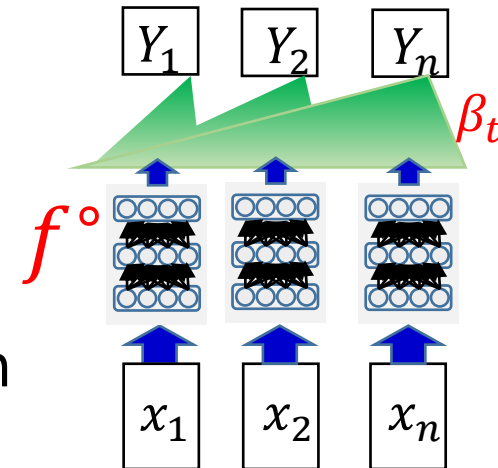
$$F_t^\circ(x) = \textcolor{red}{\beta_t^\top} \, \textcolor{blue}{f^\circ(x)}$$

- **Pretraining: Learning feature map** $\textcolor{blue}{[f^\circ]}$
  - ➢ Fourier basis, B-Spline
  - ➢ Independent of context ($t$)
  - ➢ Obtain the most "efficient" basis to represent data
    - → Internal layers

> - **Good representation**
> - **Distribution of** $\beta_t$

- **In-context learning: Estimating coefficient** $\textcolor{red}{[\beta_t]}$
  - ➢ Dependent on context ($t$)
  - ➢ Estimate the context $\beta_t$ from the instruction (Attention)
    - → Attention layer

$Y_1$  $Y_2$  $Y_n$

$\beta_t$

$f^\circ$

$x_1$  $x_2$  $x_n$

✓ Guo et al. 2023 and von Oswald et al. 2023 observed that real Transformers extract nonlinear features at lower layers and perform linear regression deeper layers.
  → It is not like performing gradient descent at every layer as in Bai et al. 2023.

## A. Nonlinear feature map (FNN)

We approximate the infinite dimensional nonlinear feature map $f°$ by DNN:

$$\phi : \mathbb{R}^d \to \mathbb{R}^N$$
$$(f° \simeq \phi)$$

Deep neural network (nonlinear feature map)

## B-1. Soft-max attention model

$$\sum_{i=1}^{n} \overset{\text{Value}}{y_{i,t}} \frac{\exp(\phi(x_{i,t})^{\top} KQ \phi(x_{\mathrm{qr},t}))}{\sum_{i'=1}^{n} \exp(\phi(x_{i',t})^{\top} KQ \phi(x_{\mathrm{qr},t}))}$$
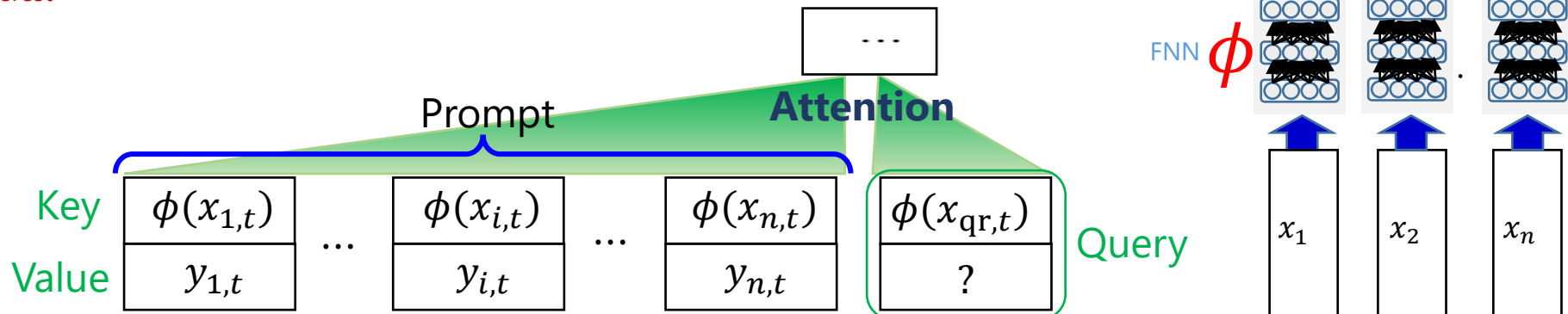
Key       Query

Predict → $y_{\mathrm{qr},t}$

## B-2. Linear attention model

[Ahn et al.: Linear attention is (maybe) all you need (to understand transformer optimization). arXiv:2310.01082]

$$\frac{1}{n} \sum_{i=1}^{n} y_{i,t} \phi(x_{i,t})^{\top} \boxed{KQ} \phi(x_{\mathrm{qr},t})$$

$\Gamma$

★ Today's interest



Attention

$\Gamma$

FNN $\phi$

Prompt       Attention

Key   $\phi(x_{1,t})$ ... $\phi(x_{i,t})$ ... $\phi(x_{n,t})$   $\phi(x_{\mathrm{qr},t})$   Query
Value   $y_{1,t}$       $y_{i,t}$       $y_{n,t}$       ?

※ In practice, each token should be a couple $(\phi(x), y)$. But, for this theoretical research, we simplify the $Q, K, V$ to a specific form

**(Linear) attention can implement linear regression:**

$$\overbrace{Y^\top \phi(X)(\phi(X)^\top \phi(X) + n\Lambda)^{-1}\phi(x_{\mathrm{qr}}) = \frac{1}{n}\sum_{i=1}^{n} y_i \phi(x_i)^\top \underbrace{\left(\frac{\phi(X)^\top \phi(X)}{n} + \Lambda\right)^{-1}}\phi(x_{\mathrm{qr}})}^{\beta^\top}}$$

$$\simeq \Gamma \text{ (prior information)}$$

Carefully chosen $\Gamma$ yields (nearly) Bayes optimal estimator.

[Gang et al. 2022; Akyurek et al. 2023; Zhang et al. 2023; Ahn et al. 2023; Mahankali et al., 2023; Wu et al. 2024]

**Empirical ICL risk** :

$$\widehat{\mathcal{L}}(\phi, \Gamma) := \frac{1}{T}\sum_{t=1}^{T}\left(y_{\mathrm{qr},t} - \frac{1}{n}\sum_{i=1}^{n} y_{i,t}\phi(x_{i,t})^\top \Gamma \phi(x_{\mathrm{qr},t})\right)^2$$

→ **Minimize with respect to $\phi$ (feature map) and $\Gamma$ (attention param).**

**The expected ICL risk**:

$$\mathcal{L}(\phi, \Gamma) := \qquad\qquad\qquad \left. qr)\right)^2\right]$$

(where $E$

**Question :**
- Can we obtain "optimal" expected risk?
- What is the benefit of ICL?

**Empirical risk minimizer:**

$$\min_{\Gamma \in \mathbb{R}^{N \times N}, \phi \in \text{DNN}} \widehat{\mathcal{L}}(\phi, \Gamma) := \frac{1}{T} \sum_{t=1}^{T} \left( y_{\text{qr},t} - \frac{1}{n} \sum_{i=1}^{n} y_{i,t} \phi(x_{i,t})^{\top} \Gamma \phi(x_{\text{qr},t}) \right)^2$$

$$\mathcal{F}_N := \left\{ \phi : \mathbb{R}^d \to \mathbb{R}^N \mid \phi \in \text{DNN} \right\}$$

**Empirical risk minimizer:**

$$(\hat{\phi}, \hat{\Gamma}) \leftarrow \underset{\Gamma \in \mathbb{R}^{N \times N}, \phi \in \mathrm{DNN}}{\arg\min} \widehat{\mathcal{L}}(\phi, \Gamma) := \frac{1}{T} \sum_{t=1}^{T} \left( y_{\mathrm{qr},t} - \frac{1}{n} \sum_{i=1}^{n} y_{i,t} \phi(x_{i,t})^{\top} \Gamma \phi(x_{\mathrm{qr},t}) \right)^2$$

$$\mathcal{F}_N := \{ \phi : \mathbb{R}^d \to \mathbb{R}^N \mid \phi \in \mathrm{DNN} \text{ with presprcified hyper-param} \}$$

**Assumption (informal)**

1. $\mathbb{E}[\beta_{t,j}^2] \lesssim j^{-2s-1-\epsilon}$    (Complexity of function space)

2. $\inf_{\phi \in \mathcal{F}_N} \max_{1 \le j \le N} \|f_j^\circ - \phi_j\|_\infty \lesssim \delta_N$    (Approx. error of each basis)

3. $\| \sum_{j=1}^{k} (f_j^\circ)^2 \|_\infty \lesssim k^{2r}$    (Bases are bounded)

4. $\left( f_j^\circ \right)_{j=1}^{\infty}$ are "near" orthonormal    (Bases are almost orthogonal to each other)

**Thm. (ICL risk bound; Kim, Nakamaki, TS, NeurIPS2024)**

$$\mathbb{E}[\mathcal{L}(\hat{\phi}, \hat{\Gamma})] \lesssim N^{-2s} + N^2 \delta_N^4 + N^{2r+1} \delta_N^2$$    Feature approximation error

$$+ \frac{N}{n} + \frac{N^{2r}}{n} \log(N) + \frac{N^{4r}}{n^2} \log^2(N)$$    In-context generalization gap

$$+ \frac{1}{T} \left( N^2 \log(\epsilon^{-1}) + \log(\mathcal{N}(\tfrac{\epsilon}{\sqrt{N}}, \mathcal{F}_N, \|\cdot\|_\infty)) \right) + \epsilon$$

Covering number of DNN

Pretraining generalization to estimate basis functions

- **Example (B-spline basis; $f_j^\circ$ is B-spline→Besov/Sobolev space):**

Estimator 1:
$$\mathbb{E}[\mathcal{L}(\hat{\phi}, \hat{\Gamma})] \lesssim N^{-2s} + \frac{N \log(N)}{n} + \frac{N^2 \log(N)}{T}$$

Bias-variance trade-off

$$\mathbb{E}[\mathcal{L}(\hat{\phi}, \hat{\Gamma})] \lesssim n^{-\frac{2s}{2s+1}} + \frac{n^{\frac{2}{2s+1}}}{T} \wedge T^{-\frac{2s}{2(s+1)}} + \frac{T^{\frac{2}{2(s+1)}}}{n}$$

**Minimax optimal w.r.t. $n$ (if $T$ is large)**

**Small $T$: memorization**
**Large $T$: generalization**

- **Example (Holder class basis; $f_j^\circ \in H^{\alpha'}(\mathbb{R}^d)$):**

Estimator 2 ($\Gamma$ is restricted to a diagonal matrix):
$$\mathbb{E}[\mathcal{L}(\hat{\phi}, \hat{\Gamma})] \lesssim N^{-2s} + \frac{N \log(N)}{n} + \frac{N^{1+\frac{d}{\alpha'}(1+s)} \log(N)}{T}$$

$$\mathbb{E}[\mathcal{L}(\hat{\phi}, \hat{\Gamma})] \lesssim n^{-\frac{2s}{2s+1}} + n^{\frac{1+\frac{d}{\alpha'}(1+s)}{2s+1}} T^{-1}$$

If there is no-pretraining, the minimax lower bound is

With many pretraining data, the pretrained model can outperform direct estimator.

$$\mathbb{E}[\mathcal{L}(\hat{\phi}, \hat{\Gamma})] \gtrsim \max\left\{ n^{-\frac{2s}{2s+1}}, n^{-\frac{2\alpha'}{2\alpha'+d}} \right\}$$

**Pretraining improves the error by estimating the bases in the pretraining phase**

$$\mathcal{L}(\hat{f}) := \mathbb{E}_{\beta, x_{\mathrm{qr}}} \left[ \left( F_\beta^\circ(x_{\mathrm{qr}}) - \hat{f}(x_{\mathrm{qr}}) \right)^2 \right]$$

$\hat{f}$: depending on the pretraining data $(x_{t,i}, y_{t,i})_{t=1, i=1}^{T,n}$ and new task data $(x_{T+1,i}, y_{T+1,i})_{i=1}^{n}$.

**Minimax risk:** $\inf\limits_{\hat{f}} \sup\limits_{f^\circ \in \mathcal{F}^\circ} \mathbb{E}[\mathcal{L}(\hat{f})]$

**Information theoretic lower bound:**

**Prop. (ICL risk lower bound)**

$$\inf_{\hat{f}} \sup_{f^\circ \in \mathcal{F}^\circ} \mathbb{E}[\mathcal{L}(\hat{f})] \gtrsim \delta_{n,T}^2$$

$$\geq \epsilon_{1,n}^2 + \epsilon_{2,n}^2$$

where

$$\begin{cases} \epsilon_{1,n}^2 \simeq \dfrac{V(\epsilon_{1,n}, \mathcal{F}^\circ)}{nT} \\[2mm] \epsilon_{2,n}^2 \simeq \dfrac{\epsilon_{2,n}^{-1/s}}{n} \end{cases}$$

(log-covering number) Complexity to estimate

: **the basis function $f^\circ$**

: **coefficient $\beta_T$**

We consider $f^\circ$ as a random variable "uniformly" distributed on a model:

$$\inf_{\hat{f}} \sup_{f^\circ \in \mathcal{F}^\circ} \mathbb{E}[\mathcal{L}(\hat{f})] \gtrsim \delta^2 \left( 1 - \frac{I(D_{1:T+1} \| (f^\circ, \beta_{T+1})) + \log(2)}{\log(\mathcal{N}(\delta, \{F_\beta^\circ\}))} \right)$$

**Optimal rate when the basis is known.**

**Complexity to estimate the basis**

(log-covering number)

$$\inf_{\hat{f}} \sup_{f^\circ \in \mathcal{F}^\circ} \mathbb{E}[\mathcal{L}(\hat{f})] \gtrsim n^{-\frac{2s}{2s+1}} + \frac{V(\epsilon_{1,n}, \mathcal{F}^\circ)}{nT}$$

where $\epsilon_{1,n}^2 \simeq \dfrac{V(\epsilon_{1,n}, \mathcal{F}^\circ)}{nT}$

- **Basis functions in Holder space** ($f_j^\circ \in H^{\alpha'}(\mathbb{R}^d)$): $\quad \dfrac{V(\epsilon_{1,n}, \mathcal{F}^\circ)}{nT} \simeq \dfrac{\epsilon_{1,n}^{-d/\alpha'}}{nT}$

$$\boxed{\inf_{\hat{f}} \sup_{f^\circ \in \mathcal{F}^\circ} \mathbb{E}[\mathcal{L}(\hat{f})] \gtrsim n^{-\frac{2s}{2s+1}} + (nT)^{-\frac{2\alpha'}{2\alpha'+d}}}$$

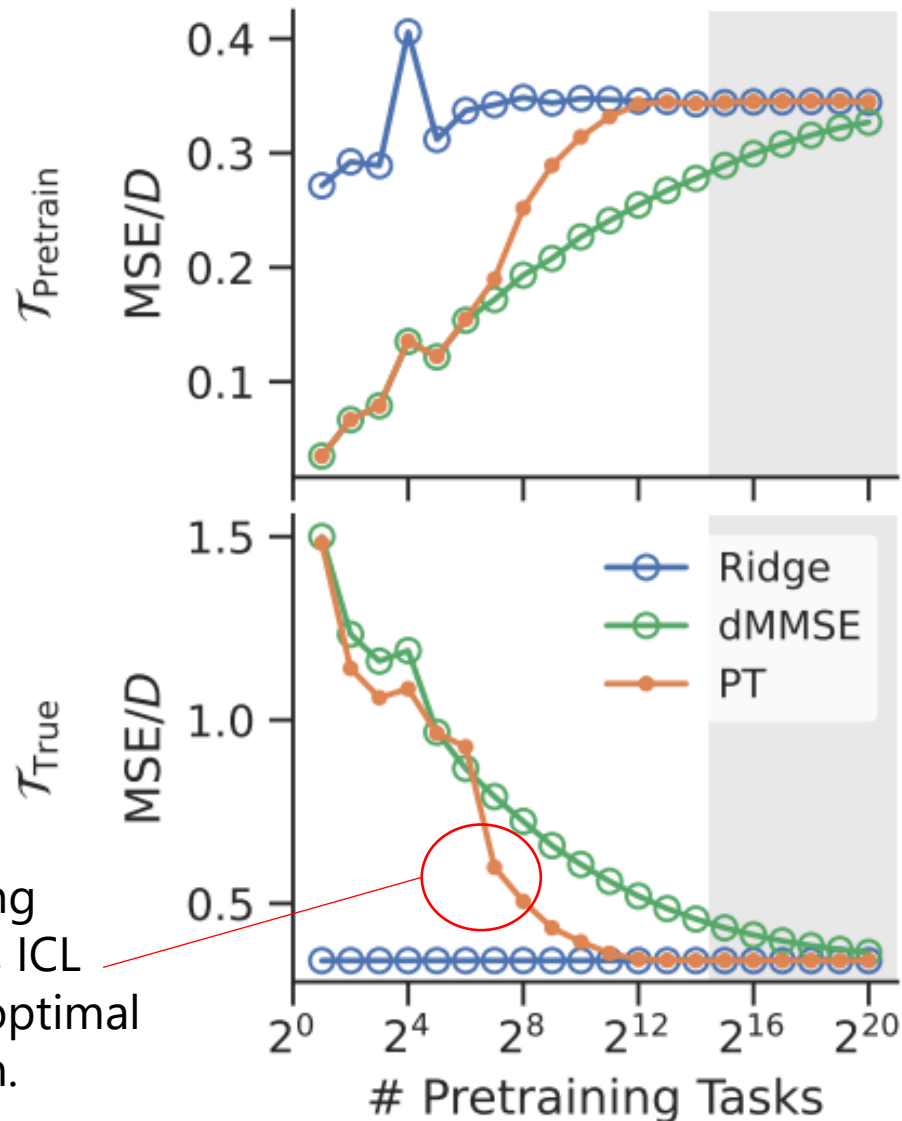Suppose that $\alpha'/d < s$, then

No pretraining ($T = 1$): $\quad n^{-\frac{2\alpha'}{2\alpha'+d}}$ 👎

$\vee$

Pretraining setting ($T \gg 1$): $\quad n^{-\frac{2s}{2s+1}}$ 👍

**When $T$ is large, pretraining can give better generalization for test instruction than learning from scratch**

# Task diversity matters



If # of pretraining tasks is enough, ICL coincides with optimal ridge regression.

[Raventós, Paul, Chen, Ganguli: Pretraining task diversity and the emergence of non-Bayesian in-context learning for regression. 2023 ]

# Presentation overview

**Statistics**

**Minimax optimality**
- Nonparametric analysis
- Approximation error analysis

**Optimization**

**Global optimality of nonlinear feature learning**
- Mean field limit
- Strict saddle

**Statistics/Optimization**

**Feature learning with one step GD**
- Single index model
- Information exponent
- Advantage of pre-training

- [Minimax optimality and approximation error bound] Kim, Nakamaki, Suzuki: Transformers are Minimax Optimal Nonparametric In-Context Learners. NeurIPS2024
- [Optimization in mean field limit] Kim, Suzuki: Transformers Learn Nonlinear Features In Context: Nonconvex Mean-field Dynamics on the Attention Landscape. ICML2024 (arXiv:2402.01258).
- [Identifying low dimensional subspace with information exponent k] Oko, Song, Suzuki, Wu: Transformer efficiently learns low-dimensional functions in context. NeurIPS2024.

So far, we have considered approximation theory.
From now on, we discuss optimization theory.

# Global optimality of GD
# for in-context learning

[Kim, Suzuki: Transformers Learn Nonlinear Features In Context: Nonconvex Mean-field Dynamics on the Attention Landscape. **ICML2024, oral presentation** (arXiv:2402.01258)]
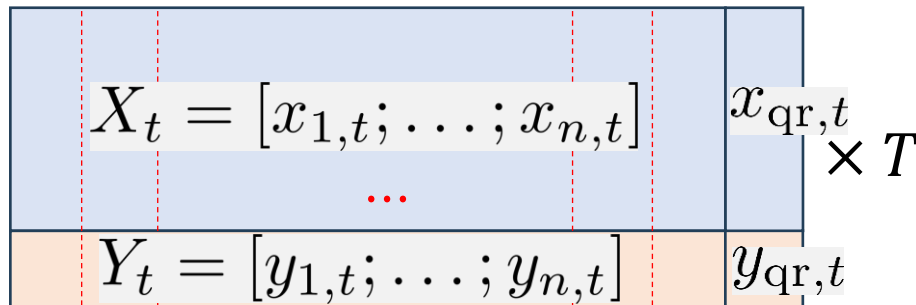


**Juno Kim**

**Model:** $\quad y_{i,t} = F_t^\circ(x_{i,t}) + \epsilon_{i,t} \qquad (i = 1, \ldots, n)$
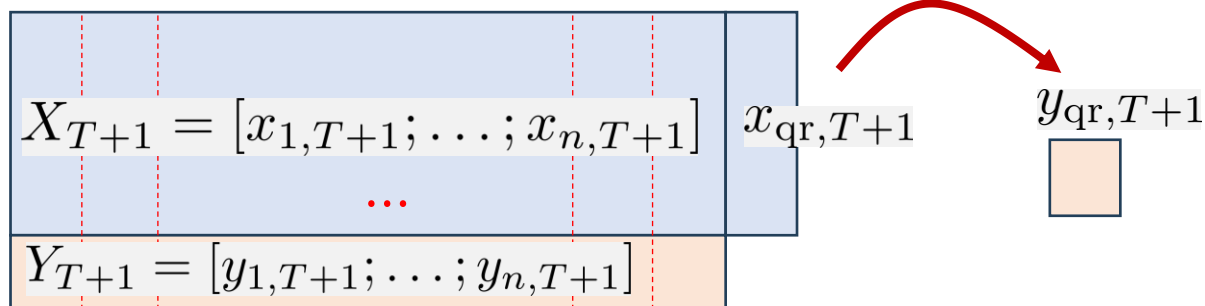
$t = 1, \ldots, T$: Task index

- The true functions $F_t$ are different across different tasks.
- $F_t^\circ$ is generated randomly for each task.

**Pretraining ($T$ tasks) :**

$$X_t = [x_{1,t}; \ldots; x_{n,t}] \quad x_{\mathrm{qr},t}$$
$$\cdots \qquad\qquad \times T$$
$$Y_t = [y_{1,t}; \ldots; y_{n,t}] \quad y_{\mathrm{qr},t}$$

➢ We observe pretraining task data $T$ times.
➢ Each task has $n$ data.

**Test task (In-context learning) :**

**Predict**

$$X_{T+1} = [x_{1,T+1}; \ldots; x_{n,T+1}] \quad x_{\mathrm{qr},T+1}$$
$$\cdots$$
$$Y_{T+1} = [y_{1,T+1}; \ldots; y_{n,T+1}]$$

$$y_{\mathrm{qr},T+1}$$

Linear model with nonlinear features:

$$F_t^\circ(x) = v_t^\top f^\circ(x) \qquad \text{where } v_t \sim N(0, I) \text{ and } f^\circ(x) \in \mathbb{R}^k.$$

We want to estimate the nonlinear feature $f^\circ$ by pretraining.

- **Mean field neural network (Barron class):**

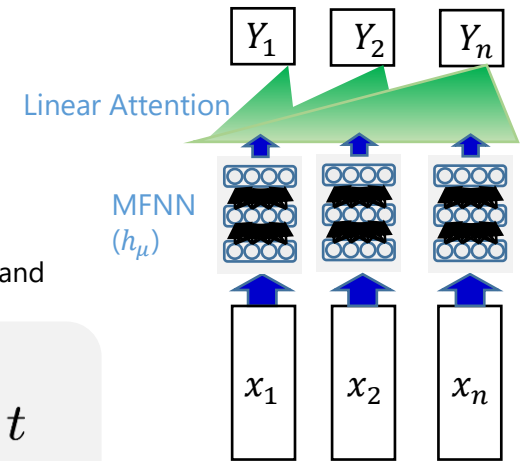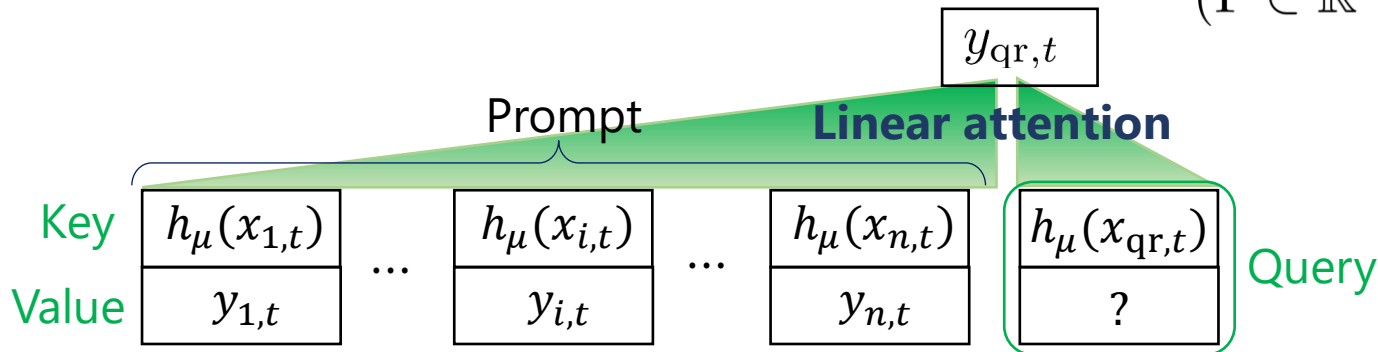$$h_\mu(x) = \int h_\theta(x) \mathrm{d}\mu(\theta) \in \mathbb{R}^k$$

$$h_\theta(x) = \mathbf{a}\sigma(\mathbf{w}^\top x) \quad (\theta = (\mathbf{a}, \mathbf{w}) \in \mathbb{R}^k \times \mathbb{R}^d)$$

- **Linear attention:** [Ahn et al.: Linear attention is (maybe) all you need (to understand transformer optimization). arXiv:2310.01082]

$$\frac{1}{n} \sum_{i=1}^{n} y_{i,t} h_\mu(x_{i,t})^\top \Gamma h_\mu(x_{\mathrm{qr},t})$$

Value    Key        Query     **Predict** $\longrightarrow$   $y_{\mathrm{qr},t}$

$$(\Gamma \in \mathbb{R}^{k \times k})$$



Linear Attention

MFNN ($h_\mu$)

$Y_1$   $Y_2$   $Y_n$

$x_1$   $x_2$   $x_n$

$y_{\mathrm{qr},t}$

Prompt    **Linear attention**

Key   | $h_\mu(x_{1,t})$ | ... | $h_\mu(x_{i,t})$ | ... | $h_\mu(x_{n,t})$ | $h_\mu(x_{\mathrm{qr},t})$ | Query

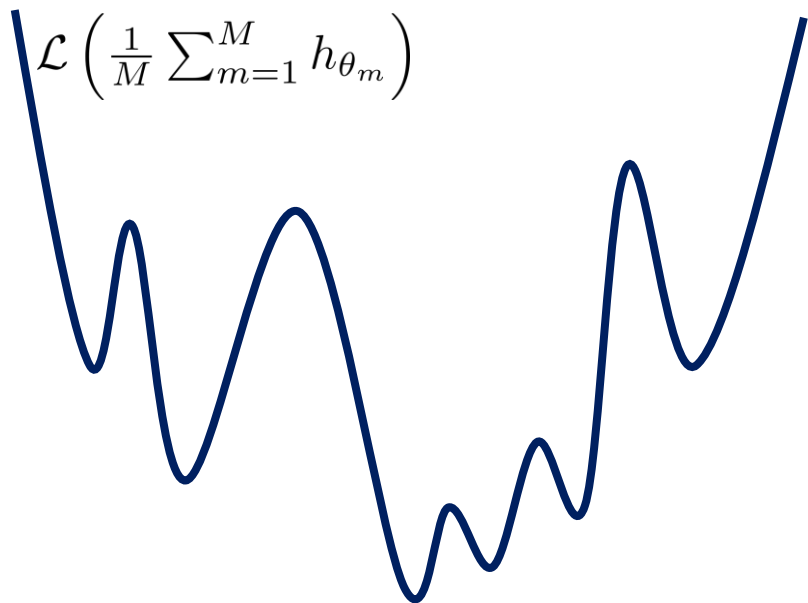Value   | $y_{1,t}$ | | $y_{i,t}$ | | $y_{n,t}$ | ? |

$$\frac{1}{M} \sum_{m=1}^{M} h_{\theta_m}(x) \rightarrow \int h_\theta(x) \mathrm{d}\mu(\theta)$$
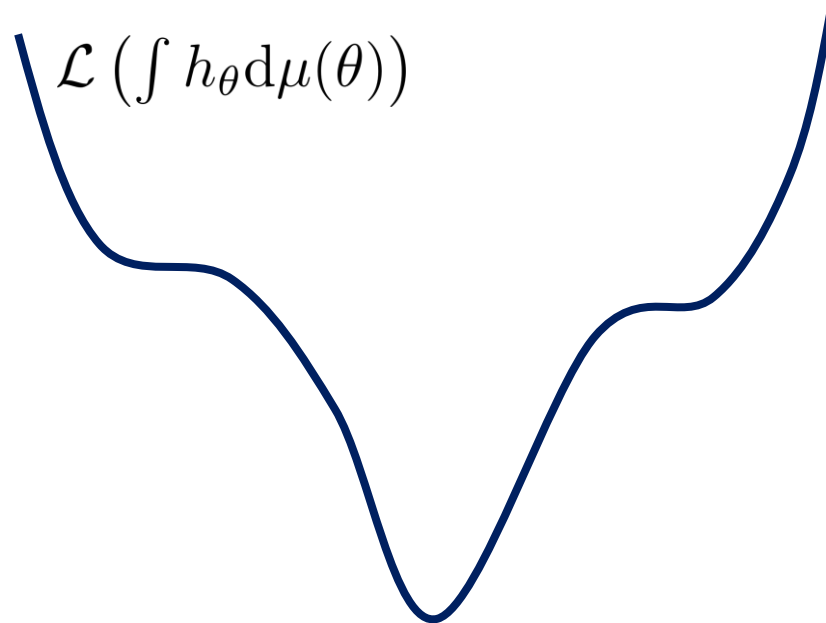
(Non-linear w.r.t. $(\theta_m)_{m=1}^M$)          (Linear w.r.t. $\mu$)

$\mathcal{L}\left(\frac{1}{M}\sum_{m=1}^M h_{\theta_m}\right)$

$\mathcal{L}\left(\int h_\theta \mathrm{d}\mu(\theta)\right)$

As a function of $\theta$          As a function of $\mu$

- Mean field Langevin dynamics: [Nitanda,Wu,Suzuki, 2022; Chizat, 2022]
  $\rightarrow$ Linear convergence with a log-Sobolev inequality for optimizing 2-layer NN.

$$\mathcal{L}(\mu_t) - \mathcal{L}^* \leq \exp(-\lambda\alpha t)(\mathcal{L}(\mu_0) - \mathcal{L}^*)$$

**Empirical ICL risk** :

$$\widehat{\mathcal{L}}(\mu, \Gamma) := \frac{1}{T} \sum_{t=1}^{T} \left( y_{\mathrm{qr},t} - \frac{1}{n} \sum_{i=1}^{n} y_{i,t} h_\mu(x_{i,t})^\top \Gamma h_\mu(x_{\mathrm{qr},t}) \right)^2$$

→ Minimize with respect to $\mu, \Gamma$.

**The expected ICL risk**:   (Large sample limit: $n \to \infty$ and $T \to \infty$)

$$\mathcal{L}(\mu, \Gamma) := \mathbb{E}_{x_{\mathrm{qr}}} \left[ \left\| f^\circ(x_{\mathrm{qr}}) - \mathbb{E}_x[f^\circ(x) h_\mu(x)^\top] \Gamma h_\mu(x_{qr}) \right\|^2 \right]$$

(note that $y_{i,t} = v_t^\top f^\circ(x_{i,t})$)

**Question :** Can we optimize $\mu, \Gamma$ by a gradient descent?
(Infinite-dimensional non-convex problem)

There have been many work on optimization guarantee on ICL for **linear model**: Zhang et al., (2023), Mahankali et al. (2023), Guo et al. (2023) to name a few.
Bu, this is a **nonlinear feature learning**.

Feature covariance $\quad \Sigma_{\mu,\nu} := \mathbb{E}_X[h_\mu(X)h_\nu^\top(X)]$

$$h_\mu(x) := \int h_\theta(x)\mathrm{d}\mu(\theta)$$

**Assumption (realizability of the true feature)**

There exists $\mu^\circ$ such that $f^\circ = h_{\mu^\circ}$ and $\Sigma_{\mu^\circ,\mu^\circ} \propto I_k$.

**Two time-scale dynamics ($\Gamma$ is optimized first):**

$$\mathcal{L}(\mu) := \min_\Gamma \mathcal{L}(\mu, \Gamma) = \min_\Gamma \mathbb{E}_{x_{\mathrm{qr}}}\left[\left\|f^\circ(x_{\mathrm{qr}}) - \mathbb{E}_x[f^\circ(x)h_\mu(x)^\top]\Gamma h_\mu(x_{qr})\right\|^2\right]$$

$$= \mathbb{E}_{x_{\mathrm{qr}}}\left[\left\|f^\circ(x_{\mathrm{qr}}) - \Sigma_{\mu^\circ,\mu}\Sigma_{\mu,\mu}^{-1}h_\mu(x_{qr})\right\|^2\right]$$

- $\mu$ is the minimizer iff $h_\mu = R h_{\mu^\circ}$ for an invertible matrix $R$
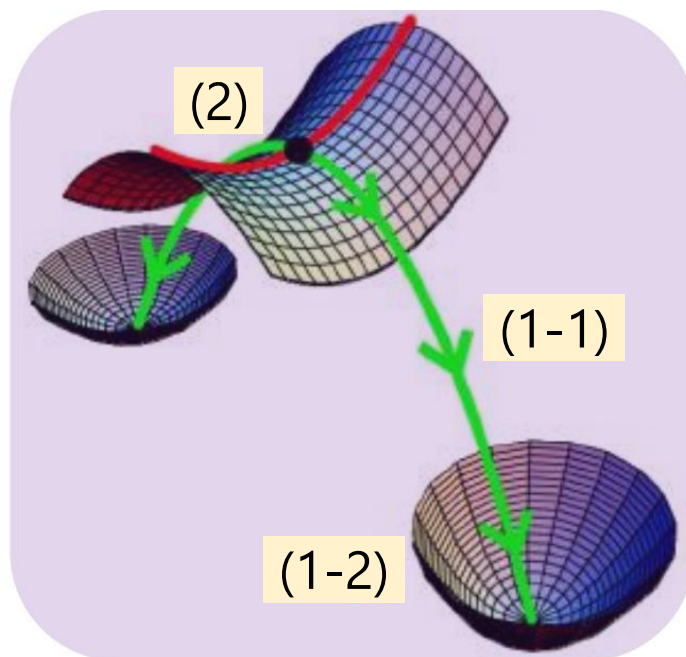
**Wasserstein gradient flow to minimize $\mathcal{L}$:**

- $\partial_t \mu_t = \nabla \cdot \left(\mu_t \nabla \dfrac{\delta\mathcal{L}(\mu_t)}{\delta\mu}\right)$

- $\dfrac{\mathrm{d}\theta_t}{\mathrm{d}t} = -\nabla\dfrac{\delta\mathcal{L}(\mu_t)}{\delta\mu}(\theta_t) \quad (\mu_t = \mathrm{Law}(\theta_t))$

- There is no spurious local minima.
- All critical points are saddle and have negative curvature.

**Theorem 1 (Strict saddle property of the loss landscape)**



There exists a **descent direction** or **negative curvature**.
Analogous to matrix completion [Ge et al., 2016, 2017; Bhojanapalli et al. 2016; Li et al., 2019].

For an orthogonal matrix $\mathbf{R} \in O(k)$, define $R\#\mu$ as the push-forward of $\mu$ along the rotation $\mathbf{R}: (a, w) \mapsto (\mathbf{R}a, w)$, i.e., $h_{\mathbf{R}\#\mu} = \mathbf{R}h_\mu$.

### Theorem 1 (**Strict saddle** property of the loss landscape)

If $\mu \in \mathcal{P}$ is not the global minimum, then one of the followings holds:

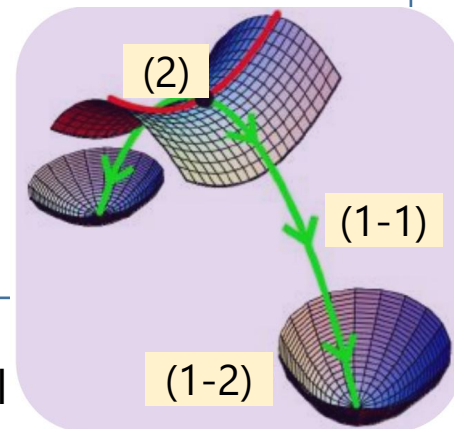**(1)** (1-1) There exists $\mathbf{R} \in \mathrm{conv}(O(k))$ such that

$$\frac{\mathrm{d}}{\mathrm{d}s}\mathcal{L}(\bar{\mu}_s)\Big|_{s=0} < 0 \quad \text{where } \bar{\mu}_s = (1-s)\mu + s\mathbf{R}\sharp\mu^\circ.$$

(1-2) Furthermore, if $0 < \mathcal{L}(\mu) < r^\circ/2$, then

$$\frac{\mathrm{d}}{\mathrm{d}s}\mathcal{L}(\bar{\mu}_s)\Big|_{s=0} \leq -\frac{4}{\|\sigma\|_\infty^2}\mathcal{L}(\mu)\left(\frac{r_0}{2} - \mathcal{L}(\mu)\right)$$

**(2)** Otherwise,

$$\mathcal{L}(\mu) > \frac{r_0}{2} \quad \text{and} \quad \frac{\mathrm{d}^2\mathcal{L}(\bar{\mu}_s)}{\mathrm{d}s^2}\Big|_{s=0} \leq -\frac{4}{k\|\sigma\|_\infty^2}\mathcal{L}(\mu)^2.$$



There exists a **descent direction** or **negative curvature**.
Analogous to matrix completion [Ge et al., 2016, 2017; Bhojanapal 2019].

Let the "Hessian" at $\mu$ be

$$H_\mu(\theta, \theta') := \nabla_\theta \nabla_{\theta'} \frac{\delta^2 \mathcal{L}(\mu)}{\delta \mu^2}(\theta, \theta')$$
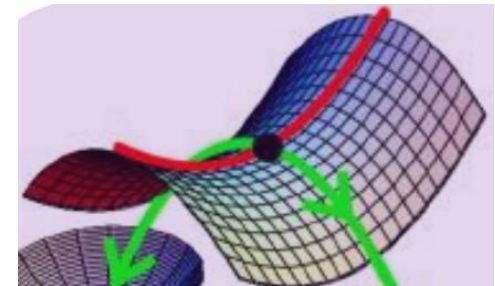
## Lemma

The Wasserstein GF $\mu_t$ around a critical point $\mu^+$ can be written as $(\text{id} + \epsilon v_t)\#\mu^+$ where the velocity field $v_t$ follows

$$\partial_t v_t(\theta) = -\int H_{\mu^+}(\theta, \theta') v_t(\theta') \mathrm{d}\mu^+(\theta') + O(\epsilon)$$

(c.f., Otto calculus)

➡ Negative curvature direction exponentially grows up!

➡ $\mu_t$ moves away from the critical point.



## Theorem (Informal)

The solution is not captured by any critical point *almost surely*.
(The solution converges to the global optimal solution almost surely)

# Decay speed of objective

Suppose that $\left\|\frac{d\mu^\circ}{d\mu_t}\right\|_\infty \le R$ (which could be ensured by using birth-death process).

**Theorem (GF moves toward a descent direction (1))**

$$\frac{d}{ds}\mathcal{L}(\bar{\mu}_s)\Big|_{s=0} < -\delta \quad \Rightarrow \quad \frac{d}{dt}\mathcal{L}(\mu_t) \le -R^{-1}\delta^2.$$

**Theorem (Accelerated convergence phase (2))**

Once $\mathcal{L}(\mu_t) \le \frac{r^\circ}{2} - \epsilon$ is satisfied,
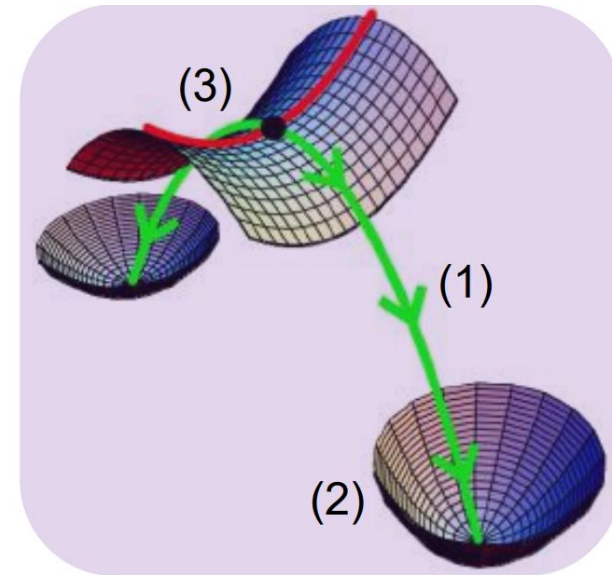$$\mathcal{L}(\mu_{t+T}) \le O\left(\frac{Rk^2}{T}\right)$$

**Theorem (Negative curvature around a saddle point (3))**

$$\frac{d^2\mathcal{L}(\bar{\mu}_s)}{ds^2} \le -\Lambda \quad \Rightarrow \quad \text{min-eigen-value}(H_{\mu_t}) \le -\Lambda/R$$
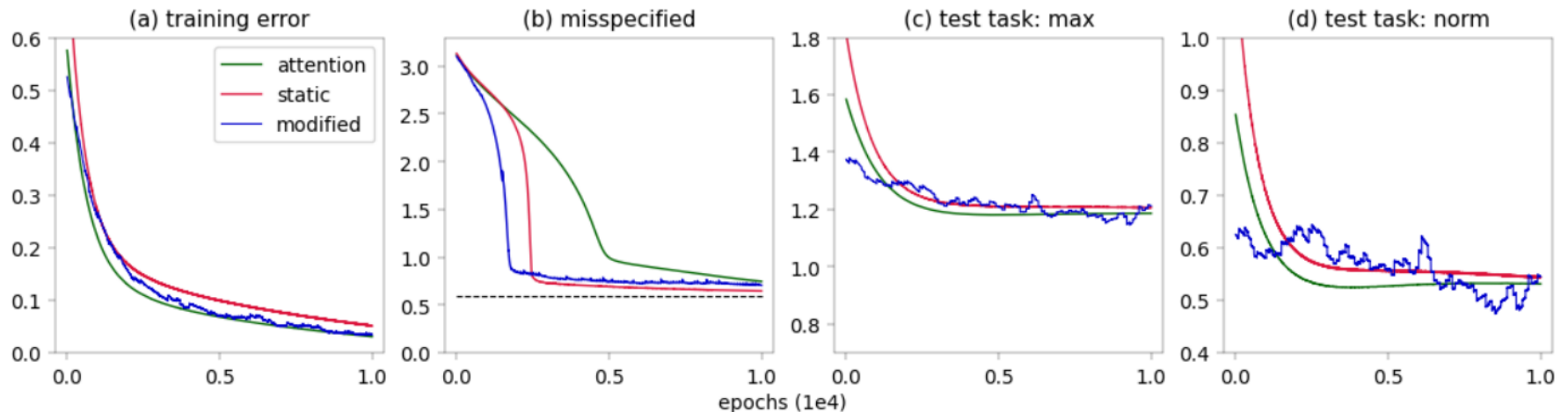
➡ Escape from the critical point exponentially fast.

We compare 3 models with $d = 20, k = 5$, and 500 neurons with sigmoid act. All models are pre-trained using SGD on 10K prompts of 1K token pairs.

1. **attention**: jointly optimizes $\mathcal{L}(\mu, \Gamma)$.
2. **static**: directly minimizes $\mathcal{L}(\mu)$.
3. **modified**: static model implementing birth-death & GP



(a) training error   (b) misspecified   (c) test task: max   (d) test task: norm

epochs (1e4)

$\rightarrow$ verify global convergence as well as improvement for misaligned model $(k_{\text{true}} = 7)$ and nonlinear test tasks $g(x) = \max_{j \leq k} h_{\mu^\circ}(x)_j$ or $g(x) = \left\| h_{\mu^\circ}(x) \right\|^2$.

**Statistics**

**Minimax optimality**
- Nonparametric analysis
- Approximation error analysis

**Optimization**

**Global optimality of nonlinear feature learning**
- Mean field limit
- Strict saddle

**Statistics/Optimization**

**Feature learning with one step GD**
- Single index model
- Information exponent
- Advantage of pre-training

- [Minimax optimality and approximation error bound] Kim, Nakamaki, Suzuki: Transformers are Minimax Optimal Nonparametric In-Context Learners. NeurIPS2024
- [Optimization in mean field limit] Kim, Suzuki: Transformers Learn Nonlinear Features In Context: Nonconvex Mean-field Dynamics on the Attention Landscape. ICML2024 (arXiv:2402.01258).
- [Identifying low dimensional subspace with information exponent k] Oko, Song, Suzuki, Wu: Transformer efficiently learns low-dimensional functions in context. NeurIPS2024.

# Nonlinear feature learning with optimization guarantee

[Oko, Song, Suzuki, Wu: Transformer efficiently learns low-dimensional functions in context. NeurIPS2024]



**Kazusato Oko**
**(The University of Tokyo/RIKEN-AIP)**
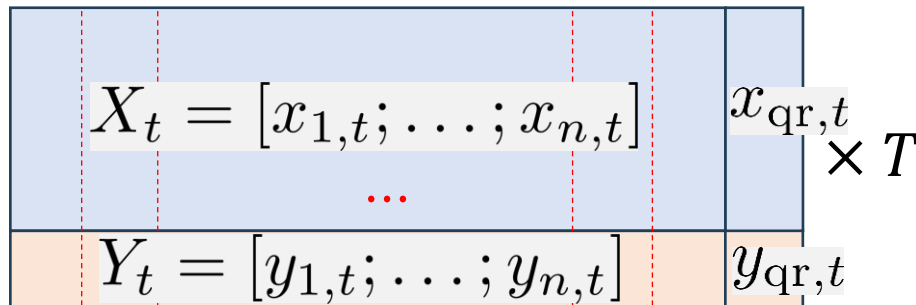


**Yujin Song**
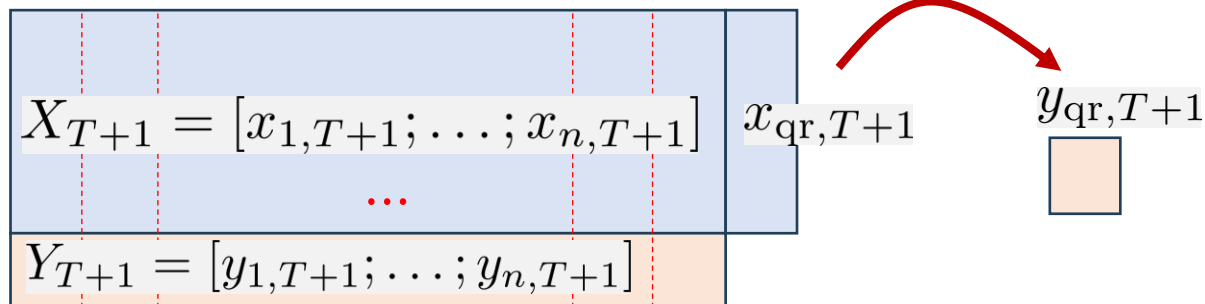**(The University of Tokyo)**



**Denny Wu**
**(NYU/Flatiron Institute)**

**Model:** $y_{i,t} = f_*^t(x_{i,t}) + \epsilon_{i,t} \qquad (i = 1, \ldots, n)$

$t = 1, \ldots, T$: Task index

**Pretraining ($T$ tasks) :**

$$X_t = [x_{1,t}; \ldots; x_{n,t}] \quad x_{\mathrm{qr},t}$$
$$\ldots \quad \times T$$
$$Y_t = [y_{1,t}; \ldots; y_{n,t}] \quad y_{\mathrm{qr},t}$$

➤ We observe pretraining task data $T$ times.
➤ Each task has $n$ data.

**Test task (In-context learning) :**

**Predict**

$$X_{T+1} = [x_{1,T+1}; \ldots; x_{n,T+1}] \quad x_{\mathrm{qr},T+1} \qquad y_{\mathrm{qr},T+1}$$
$$\ldots$$
$$Y_{T+1} = [y_{1,T+1}; \ldots; y_{n,T+1}]$$

**Gaussian single index model:**

$$f_*^t(x) = \sigma_*^t(\langle x, \beta_t \rangle)$$

where the link $\sigma_*^t$ and the direction $\beta_t$ are generated randomly:

---

**$\beta_t$**   $\beta_t$ is distributed uniformly on a unit sphere in an $r < d$ dimensional <u>linear subspace $\mathcal{S}$</u>:

$$\beta_t \sim \mathrm{Unif}(\mathrm{Unit}(\mathcal{S})) \text{ where } \dim(\mathcal{S}) = r \ll d$$

---

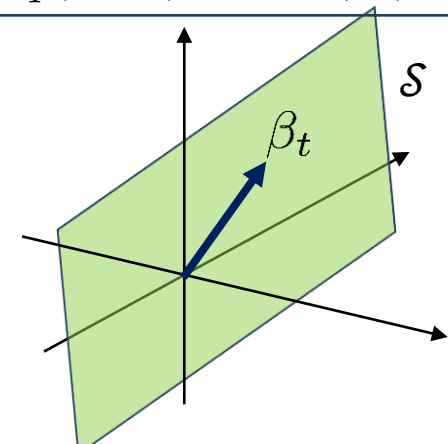**$\sigma_*^t$**   $\sigma_*^t(z) = \sum_{i=k}^P c_i^t \mathrm{He}_i(z)$

where $c_i^t$ is randomly generated from a distribution satisfying

$$\mathbb{E}[c_2^t] \neq 0, \ \sum_{i=2}^P (c_i^t)^2 = \Theta(1) \text{ (a.s.)}, \ (c_2^t, \dots, c_P^t) \neq (0, \dots, 0) \text{ (a.s.)}$$

---

$\Rightarrow$ **Information exponent = $k$.**

**The feature has a low dimensional structure.**

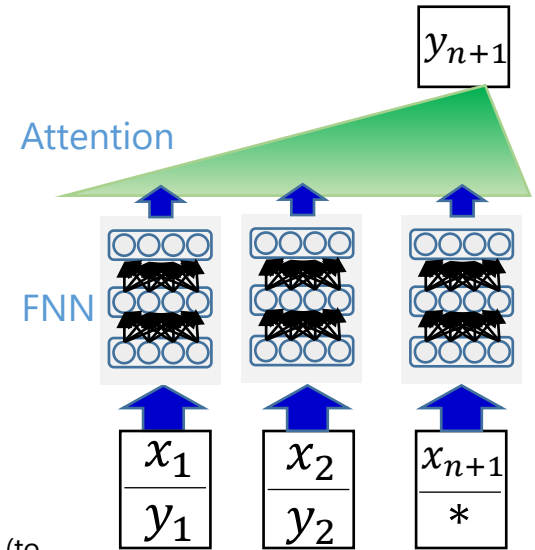We want to estimate the subspace $\mathcal{S}$ and the basis functions $\mathrm{He}_i$ in the pretraining stage.

- **FNN layer** $(f_W : \mathbb{R}^d \to \mathbb{R}^m)$ :

$$f_{\mathbf{w},\mathbf{b}}(x) = \begin{pmatrix} \sigma(\mathbf{w}_1^\top x + b_1) \\ \sigma(\mathbf{w}_2^\top x + b_2) \\ \vdots \\ \sigma(\mathbf{w}_m^\top x + b_m) \end{pmatrix} =: \sigma(\mathbf{W}x + \mathbf{b})$$
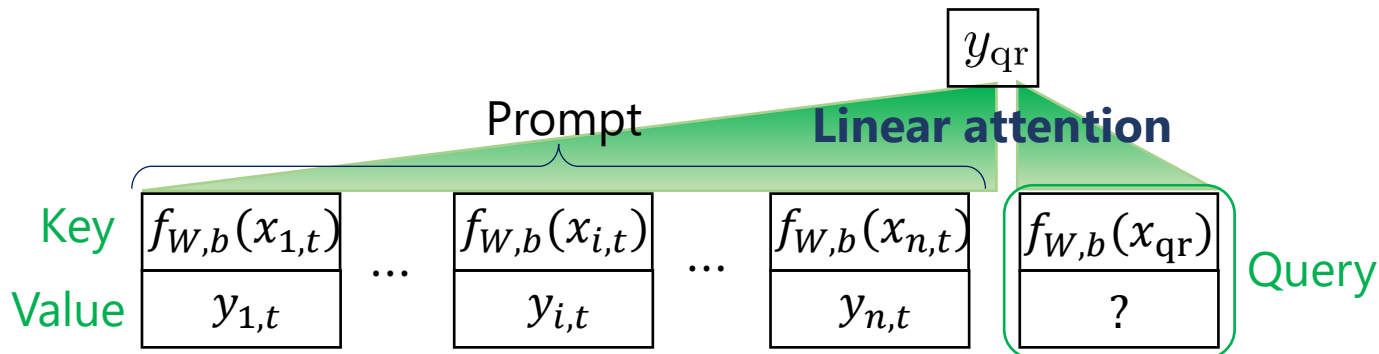
($\sigma$: ReLU)

Attention

FNN

$$\begin{array}{ccc} \dfrac{x_1}{y_1} & \dfrac{x_2}{y_2} & \dfrac{x_{n+1}}{*} \end{array}$$

$y_{n+1}$

- **Linear attention model:** [Ahn et al.: Linear attention is (maybe) all you need (to understand transformer optimization). arXiv:2310.01082]

$$f(X_t, Y_t, x; \mathbf{W}, \mathbf{b}, \Gamma) = \\ \frac{1}{n} \sum_{i=1}^{n} \underbrace{y_{i,t}}_{\text{Value}} \underbrace{f_{\mathbf{w},\mathbf{b}}(x_{i,t})^\top}_{\text{Key}} \Gamma \underbrace{f_{\mathbf{w},\mathbf{b}}(x_{\mathrm{qr}})}_{\text{Query}}$$

(linear regression)

$\longrightarrow$ **Predict** $y_{\mathrm{qr}}$ $(\Gamma \in \mathbb{R}^{k \times k})$

$y_{\mathrm{qr}}$

Prompt       **Linear attention**

| Key | $f_{W,b}(x_{1,t})$ | ... | $f_{W,b}(x_{i,t})$ | ... | $f_{W,b}(x_{n,t})$ | $f_{W,b}(x_{\mathrm{qr}})$ | Query |
| --- | --- | --- | --- | --- | --- | --- | --- |
| Value | $y_{1,t}$ | | $y_{i,t}$ | | $y_{n,t}$ | ? | |

$$E = \begin{pmatrix} \sigma(\mathbf{w}_1^\top x_1 + b_1) & \cdots & \sigma(\mathbf{w}_1^\top x_n + b_1) & \sigma(\mathbf{w}_1^\top x_{n+1} + b_1) \\ \vdots & \ddots & \vdots & \vdots \\ \sigma(\mathbf{w}_m^\top x_1 + b_m) & \cdots & \sigma(\mathbf{w}_m^\top x_n + b_m) & \sigma(\mathbf{w}_m^\top x_{n+1} + b_m) \\ y_1 & \cdots & y_n & 0 \end{pmatrix}$$



Attention

FNN

$x_1$ / $y_1$   $x_2$ / $y_2$   $x_{n+1}$ / *   $y_{n+1}$

$$f_{\mathrm{Att}}(X, Y) = W_V E \cdot \mathrm{softmax}\left(\frac{(W_K E)^\top W^Q E}{\lambda}\right)$$

$$= \frac{1}{C_{n+1}} \sum_{j=1}^{n} (W_V E_{:,j}) \exp\left(\frac{(W_K E_{:,j})^\top (W_Q E_{:,n+1})}{\lambda}\right)$$

Consider the following special setting:

$$W_V = \begin{bmatrix} 0_{1\times m} & 1 \end{bmatrix} \qquad W_K^\top W_Q = \begin{pmatrix} \Gamma & * \\ 0_{1\times m} & * \end{pmatrix}$$

Then,

$$f_{\mathrm{Att}}(X, Y) = \frac{1}{C_{n+1}} \sum_{j=1}^{n} y_j \exp\left(f_{\mathbf{w}, \mathbf{b}}(x_j)^\top \Gamma f_{\mathbf{w}, \mathbf{b}}(x_{n+1})\right)$$

By ignoring the normalization constant $C_{n+1}$ and the nonlinear term exp, we obtain the linear attention in the previous slide.

**Empirical ICL risk** :

$$\widehat{\mathcal{L}}(\mathbf{W}, \mathbf{b}, \Gamma) := \frac{1}{T} \sum_{t=1}^{T} \left( y_{\mathrm{qr},t} - \frac{1}{n} \sum_{i=1}^{n} y_{i,t} f_{\mathbf{W},\mathbf{b}}(x_{i,t})^{\top} \Gamma f_{\mathbf{W},\mathbf{b}}(x_{\mathrm{qr},t}) \right)^2$$

→ Minimize with respect to $W, b, \Gamma$.

**The expected ICL risk**:    (Large sample limit: $n \to \infty$ and $T \to \infty$)

$$\mathcal{L}(\mathbf{W}, \mathbf{b}, \Gamma) := \mathbb{E}_{x_{\mathrm{qr}}, f_*} \left[ \left( f_*(x_{\mathrm{qr}}) - \mathbb{E}_x[f_*(x) f_{\mathbf{W},\mathbf{b}}(x)^{\top}] \Gamma f_{\mathbf{W},\mathbf{b}}(x_{qr}) \right)^2 \right]$$

(note that $y_{i,t} = f_*^t(x_{i,t}) + \epsilon_{i,t}$)

**Question :**

- Can we estimate $W, b, \Gamma$ by gradient descent?   (Non-convex problem)
- How large is the sample complexity?

Initialize $\mathbf{w}_j^{(0)} \sim \text{Unif}(\mathbb{S}^{d-1})$, $b_j = 0$, $\Gamma_{j,j}^{(0)} = \text{Unif}(\{\pm 1\})$ (diagonal).

## • Stage 1: One-step gradient descent.

Optimize $W$ by a **one-step gradient descent**:

Find the subspace $\mathcal{S}$

$$\mathbf{w}_j^{(1)} \leftarrow \mathbf{w}_j^{(0)} - \eta \left[ \nabla_{\mathbf{w}_j} \frac{1}{T_1} \sum_{t=1}^{T_1} \left( y_{\mathrm{qr},t} - f(X_t, Y_t, x_{\mathrm{qr},t}; \mathbf{W}^{(0)}, \mathbf{b} = 0, \Gamma^{(0)}) \right)^2 + \lambda \mathbf{w}_j^{(0)} \right]$$

➤ Analogous to one-step GD for 2-layer NN [Damian et al. 22; Ba et al. 22].
➤ Since the true link function has $\mathrm{IE} = 2$, we can recover the subspace $\mathcal{S}$ by one-step GD with large step size.

## • Stage 2: Optimization of $\Gamma$.

Randomly re-initialize $b_j \sim \text{Unif}([-1,1])$.
Optimize $\Gamma$ based on the feature $W$ obtained at Stage 1:

$$\widehat{\Gamma} \leftarrow \arg\min_\Gamma \left\{ \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} \left( y_{\mathrm{qr},t} - f(X_t, Y_t, x_{\mathrm{qr},t}; \mathbf{W}^{(1)}, \mathbf{b}, \Gamma) \right)^2 + \lambda \|\Gamma\|_F^2 \right\}$$

$$\frac{1}{n} \sum_{i=1}^{n} y_{i,t} f_{\mathbf{W},\mathbf{b}}(x_{i,t})^\top \Gamma f_{\mathbf{W},\mathbf{b}}(x_{\mathrm{qr}})$$

Train the attention to extract the coefficient $\beta_t$.

$$\widehat{\Gamma} \leftarrow \arg\min_{\Gamma} \left\{ \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} \left( y_{\mathrm{qr},t} - f(X_t, Y_t, x_{\mathrm{qr},t}; \mathbf{W}^{(1)}, \mathbf{b}, \Gamma) \right)^2 + \lambda \|\Gamma\|_F^2 \right\}$$

$$f_t(X_t, Y_t, x_{\mathrm{qr},t}; \mathbf{W}, \mathbf{b}, \Gamma) = \frac{1}{n} \sum_{i=1}^{n} y_{i,t} f_{\mathbf{W},\mathbf{b}}(x_{i,t})^\top \Gamma f_{\mathbf{W},\mathbf{b}}(x_{\mathrm{qr},t})$$

Then, $\hat{\Gamma}$ performs the ridge regression:

$$f_t(X_t, Y_t, x_{\mathrm{qr},t}; \mathbf{W}^{(1)}, \mathbf{b}, \hat{\Gamma}) = f_{\mathbf{W}^{(1)},\mathbf{b}}(x_{\mathrm{qr},t})^\top \left( \frac{1}{nT_2} F_{T_1:T_2}^\top F_{T_1:T_2} + \lambda I \right)^{-1} F_t Y_t$$

where $F_t = [f_{\mathbf{W}^{(1)},\mathbf{b}}(x_{1,t}), \ldots, f_{\mathbf{W}^{(1)},\mathbf{b}}(x_{n,t})].$

If we can obtain **nice basis functions** $f_{W^{(1)},b}$ at Stage 1, the target function can be well estimated in the test task.

## Theorem (ICL risk bound)

Let $n^*$ be the number of examples in test task. If the one-step GD is performed with

$$T_1 = \Theta(d^{k+1}) \text{ and } n = \widetilde{\Omega}(d^k),$$

then the trained Transformer achieves the following test loss:

$$\mathcal{L}(\widehat{\mathbf{W}}, \widehat{\mathbf{b}}, \widehat{\Gamma}) \lesssim \underbrace{\frac{r^{3P/2}}{\sqrt{m}}}_{\text{Approximation error}} + \underbrace{\sqrt{\frac{r^{4P}}{T_2}}}_{\substack{\text{Error to} \\ \text{estimate } \Gamma}} + \underbrace{r^{2P}\sqrt{\frac{1}{n} + \frac{1}{n^*}}}_{\substack{\text{Error to estimate} \\ \text{in the test task}}}.$$

$m$: width of NN, $T_1$: number of tasks in Stage 1 (learning $W$), $T_2$: number of tasks in Stage 2 (learning $\Gamma$), $n$: number of examples in pretraining-task.

- Without pretraining (non-ICL setting), $n^* = \Omega(d^p)$ for kernel method and $n^* = \Omega(d^{k/2})$ for CSQ algorithm are required. But, in ICL, $n^*$ can be independent of $d$ ($n^* = \text{poly}(r)$).

- To estimate $W$, it requires $T_1 n = \Theta(d^{2k+1})$ datapoints while Damian et al. (2022) required only $\Theta(d^2)$ data points because we need enough task diversity.
  - But, ICL does not update their parameters based on the in-context examples.

## Theorem (ICL risk bound)

Let $n^*$ be the number of examples in test task. If the one-step GD is performed with

$$T_1 = \Theta(d^{k+1}) \text{ and } n = \widetilde{\Omega}(d^k),$$

then the trained Transformer achieves the following test loss:

$$\mathcal{L}(\widehat{\mathbf{W}}, \widehat{\mathbf{b}}, \widehat{\Gamma}) \lesssim \frac{r^{3P/2}}{\sqrt{m}} + \sqrt{\frac{r^{4P}}{T_2} + r^{2P}\sqrt{\frac{1}{n} + \frac{1}{n^*}}}.$$

Approximation error    Error to estimate $\Gamma$    Error to estimate in the test task

$m$:

(lea

| | w/o pretraining | | w/ pretraining |
|---|---|---|---|
| Method | Kernel | NN (CSQ or SQ) | ICL |
| Sample complexity | $d^P$ | $d^{k/2}$ (or $d$) | $r^{2P}$ |
| Pretraining | --- | --- | $T_1 = d^{k+1}, n = d^k$ |

- V
  $\Omega$
  $d$

- To esti ... al. (2022)
  require ... ty.
  ➢ Bu ... t examples.

If we observe many data during pretraining,
ICL with Transformer can generalize well in test tasks.

- The one-step GD update (with regularization) projects the initial vector $w_j^{(0)}$ to the subspace $\mathcal{S}$.
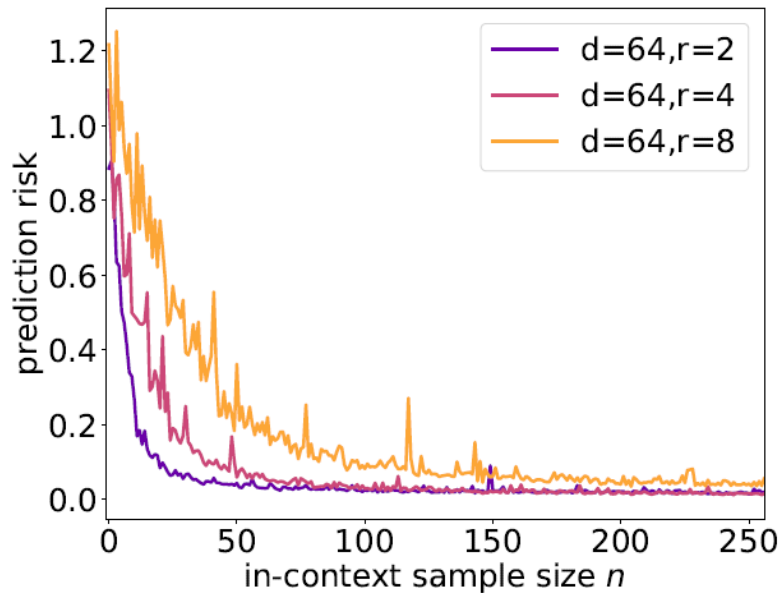


- Learning $W$ : Subspace $\mathcal{S}$ is obtained.
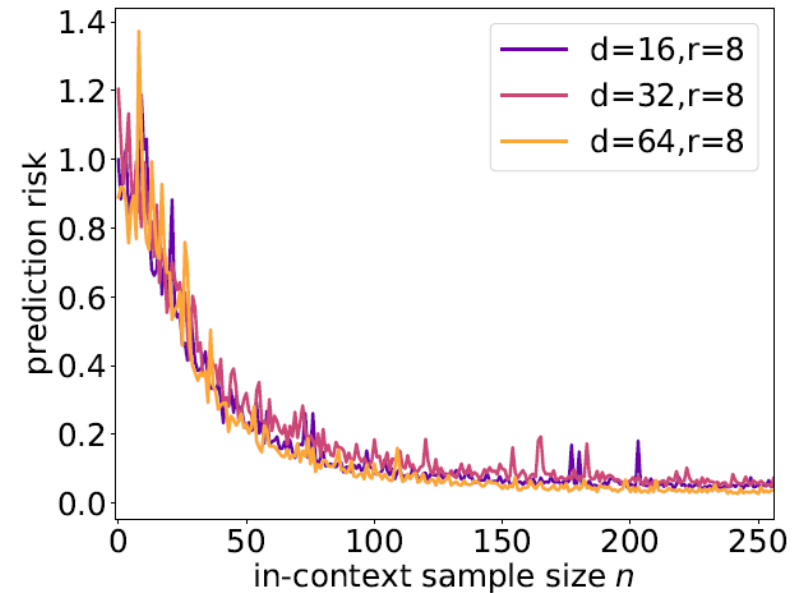- Learning $\Gamma$ : Attention to obtain the coefficients on basises.

- If we have many neurons, $\left(w_j^{(1)}\right)_{j=1}^m$ spans the subspace $\mathcal{S}$ (**1st -stage**).

- If we have sufficiently large number of neuros $\left(\sigma(w_j^{(1)\top}w + b_j)\right)_{j=1}^m$, the model can well approximate the target polynomial $\sigma^*(\langle \beta_t, x \rangle)$ by **linear combination of the ReLU-neurons (2nd-stage + test prompt)**.

Fixing d, changing r

Fixing r, changing d

GPT2 model with 12-layers (~22M parameters)

Only $r$ affects the result, $d$ does not.

- Learning theory of in-context learning

$$F_t^\circ(x) = \beta_t^\top f^\circ(x)$$

  ➢ **Pretraining:** Obtaining nonlinear feature $[f^\circ]$

  ➢ **In-context instruction:** Learning coefficient $[\beta_t]$

- Nonparametric regression theory
  ➢ Minimax optimality
  ➢ Task diversity matters.

- Optimization theory
  ➢ Feature learning by mean-field neural network
  ➢ Estimating single-index model by gradient descent
  → Feature learning helps to improve the sample complexity of in-context learning.