

Trained Transformers Learn Linear Models In-Context (2024 JMLR)

Presenter: Gao Zhuo

HSBC Business School, Peking University

2025/08/17

1 Introduction

2 Preliminaries

3 Main results

In-context Learning (ICL)

Prompt

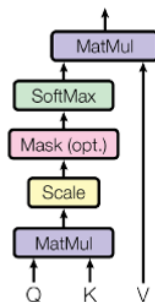
$$P = (x_1, h(x_1), \dots, x_N, h(x_N), x_{query})$$

where:

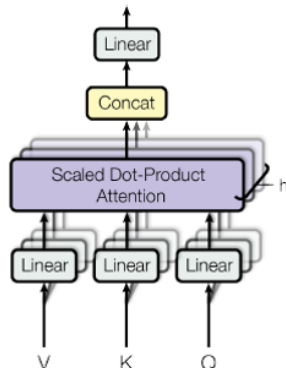
- x_i, x_{query} is sampled i.i.d. from a distribution D_x
 - h is sampled independently from a distribution over functions in a function class H
 - $h(x_i)$ is the output for each input sample data x_i in the prompt
-
- **ICL:** when given a short sequence of input-output pairs (called a prompt) from a particular task as input, the model can formulate predictions on test examples without having to make any updates to the parameters in the model.

Transformer

Scaled Dot-Product Attention



Multi-Head Attention



- The core part of transformer: self-attention module is shown above. Can we just use this module to perform ICL procedure?

Transformer

Attention + Residual Connection

$$f_{Attn} = E + W^P W^V E \cdot \text{softmax}\left(\frac{(W^K E)^T W^Q E}{\rho}\right)$$

where:

- $E \in R^{d_e, d_N}$ is an embedding matrix formed using a prompt $(x_1, y_1, \dots, x_N, y_N, x_{query})$
- One natural way to form E is to stack $(x_i, y_i) \in R^{d+1}$ as the first N columns of E and to let the final column be $(x_{query}, 0) \in R^{d+1}$
- W^K, W^Q, W^V is the key, query, and value weight matrices
- W^P is the projection matrix
- $\rho > 0$ is a normalization factor
- Softmax is applied column-wise and, given a vector input of v , the i -th entry of $\text{softmax}(v)$ is given by $\exp(v_i) / \sum_s \exp(v_s)$

In-context Learning (ICL)

Prompt

$$P = (x_1, h(x_1), \dots, x_N, h(x_N), x_{query})$$

where:

- x_i, x_{query} is sampled i.i.d. from a distribution D_x
 - h is sampled independently from a distribution over functions in a function class H
 - $h(x_i)$ is the output for each input sample data x_i in the prompt
-
- The behavior of the trained transformers can mimic those of familiar learning algorithms like ordinary least squares
 - What about the case when the mapping function h is linear?
 $h(x) = \langle w, x \rangle$

In-context Learning (ICL)

Prompt in Linear Form

$$P = (x_1, \langle w, x_1 \rangle, \dots, x_N, \langle w, x_N \rangle, x_{query})$$

where:

- x_i, x_{query} is sampled i.i.d. from a distribution D_x
 - w is weighted vector of standard Gaussian distribution
-
- Transformers based neural network architectures which are capable of achieving small prediction error for query examples
 - However, how transformer architectures produce models which are capable of in-context learning?

1 Introduction

2 Preliminaries

3 Main results

In-context Learning

Definition 1 (Trained on in-context examples) Let \mathcal{D}_x be a distribution over an input space \mathcal{X} , $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$ a set of functions $\mathcal{X} \rightarrow \mathcal{Y}$, and $\mathcal{D}_{\mathcal{H}}$ a distribution over functions in \mathcal{H} . Let $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a loss function. Let $\mathcal{S} = \cup_{n \in \mathbb{N}} \{(x_1, y_1, \dots, x_n, y_n) : x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}$ be the set of finite-length sequences of (x, y) pairs and let

$$\mathcal{F}_{\Theta} = \{f_{\theta} : \mathcal{S} \times \mathcal{X} \rightarrow \mathcal{Y}, \theta \in \Theta\}$$

be a class of functions parameterized by θ in some set Θ . For $N > 0$, we say that a model $f : \mathcal{S} \times \mathcal{X} \rightarrow \mathcal{Y}$ is trained on in-context examples of functions in \mathcal{H} under loss ℓ w.r.t. $(\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x)$ if $f = f_{\theta^*}$ where $\theta^* \in \Theta$ satisfies

$$\theta^* \in \operatorname{argmin}_{\theta \in \Theta} \mathbb{E}_{P=(x_1, h(x_1), \dots, x_N, h(x_N), x_{\text{query}})} [\ell(f_{\theta}(P), h(x_{\text{query}}))], \quad (1)$$

where $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x$ and $h \sim \mathcal{D}_{\mathcal{H}}$ are independent. We call N the length of the prompts seen during training.

- After training, we hope that the predictor $f_{\theta}(P)$ is as close as the true value $h(x_{\text{query}})$

In-context Learning

Definition 2 (In-context learning of a hypothesis class) Let \mathcal{D}_x be a distribution on an input space \mathcal{X} , $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$ a class of functions $\mathcal{X} \rightarrow \mathcal{Y}$, and $\mathcal{D}_{\mathcal{H}}$ a distribution on functions in \mathcal{H} . Let $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a loss function. Let $\mathcal{S} = \cup_{n \in \mathbb{N}} \{(x_1, y_1, \dots, x_n, y_n) : x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}$ be the set of finite-length sequences of (x, y) pairs. We say that a model $f : \mathcal{S} \times \mathcal{X} \rightarrow \mathcal{Y}$ defined on prompts of the form $P = (x_1, h(x_1), \dots, x_M, h(x_M), x_{\text{query}})$ in-context learns a hypothesis class \mathcal{H} under loss ℓ with respect to $(\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x)$ up to error $\eta \in \mathbb{R}$ if there exists a function $M_{\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x}(\varepsilon) : (0, 1) \rightarrow \mathbb{N}$ such that for every $\varepsilon \in (0, 1)$, and for every prompt P of length $M \geq M_{\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x}(\varepsilon)$,

$$\mathbb{E}_{P=(x_1, h(x_1), \dots, x_M, h(x_M), x_{\text{query}})} \left[\ell \left(f(P), h(x_{\text{query}}) \right) \right] \leq \eta + \varepsilon, \quad (2)$$

where the expectation is over the randomness in $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x$ and $h \sim \mathcal{D}_{\mathcal{H}}$.

- i.e., for any level of precision we desire, there can exist a test prompt length that meets our requirements

Single-layer Linear Self-attention

Linear Self-attention (LSA)

$$f_{LSA} = E + W^{PV}E \cdot \frac{E^T W^{KQ} E}{\rho}$$

where:

- $E \in \mathbb{R}^{d_e, d_N}$ is an embedding matrix formed using a prompt $(x_1, y_1, \dots, x_N, y_N, x_{query})$
- One natural way to form E is to stack $(x_i, y_i) \in \mathbb{R}^{d+1}$ as the first N columns of E and to let the final column be $(x_{query}, 0) \in \mathbb{R}^{d+1}$
- $\rho > 0$ is a normalization factor
- $W^{KQ} = (W^K)^T W^Q$, $W^{PV} = W^P W^V$
- **Softmax is removed**

$$E = E(P) = \begin{pmatrix} x_1 & x_2 & \cdots & x_N & x_{query} \\ y_1 & y_2 & \cdots & y_N & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (N+1)}.$$

Single-layer Linear Self-attention

by the LSA layer, actually only part of W^{PV} and W^{KQ} affect the prediction. To see how, let us denote

$$W^{PV} = \begin{pmatrix} W_{11}^{PV} & w_{12}^{PV} \\ (w_{21}^{PV})^\top & w_{22}^{PV} \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)}, \quad W^{KQ} = \begin{pmatrix} W_{11}^{KQ} & w_{12}^{KQ} \\ (w_{21}^{KQ})^\top & w_{22}^{KQ} \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)}, \quad (5)$$

where $W_{11}^{PV} \in \mathbb{R}^{d \times d}$; $w_{12}^{PV}, w_{21}^{PV} \in \mathbb{R}^d$; $w_{22}^{PV} \in \mathbb{R}$; and $W_{11}^{KQ} \in \mathbb{R}^{d \times d}$; $w_{12}^{KQ}, w_{21}^{KQ} \in \mathbb{R}^d$; $w_{22}^{KQ} \in \mathbb{R}$. Then, the prediction \hat{y}_{query} is

$$\hat{y}_{\text{query}} = \begin{pmatrix} (w_{21}^{PV})^\top & w_{22}^{PV} \end{pmatrix} \cdot \left(\frac{EE^\top}{N} \right) \begin{pmatrix} W_{11}^{KQ} \\ (w_{21}^{KQ})^\top \end{pmatrix} x_{\text{query}}, \quad (6)$$

since only the last row of W^{PV} and the first d columns of W^{KQ} affects the prediction, which means we can simply take all other entries zero in the following sections.

- Since the prediction takes only the **right-bottom** entry of the token matrix output, we can write \hat{y}_{query} into the above form

Training Procedure

In this work, we will consider the task of in-context learning linear predictors. We will assume training prompts are sampled as follows. Let Λ be a positive definite covariance matrix. Each training prompt, indexed by $\tau \in \mathbb{N}$, takes the form

$$P_\tau = (x_{\tau,1}, h_\tau(x_{\tau,1}), \dots, x_{\tau,N}, h_\tau(x_{\tau,N}), x_{\tau,\text{query}}),$$

where task weights $w_\tau \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$, inputs $x_{\tau,i}, x_{\tau,\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Lambda)$, and labels $h_\tau(x) = \langle w_\tau, x \rangle$.

Each prompt corresponds to an embedding matrix E_τ , formed using the transformation (4):

$$E_\tau := \begin{pmatrix} x_{\tau,1} & x_{\tau,2} & \cdots & x_{\tau,N} & x_{\tau,\text{query}} \\ \langle w_\tau, x_{\tau,1} \rangle & \langle w_\tau, x_{\tau,2} \rangle & \cdots & \langle w_\tau, x_{\tau,N} \rangle & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (N+1)}.$$

We denote the prediction of the LSA model on the query label in the task τ as $\hat{y}_{\tau,\text{query}}$, which is the bottom-right element of $f_{\text{LSA}}(E_\tau)$, where f_{LSA} is the linear self-attention model defined in (3). The empirical risk over B independent prompts is defined as

$$\hat{L}(\theta) = \frac{1}{2B} \sum_{\tau=1}^B \left(\hat{y}_{\tau,\text{query}} - \langle w_\tau, x_{\tau,\text{query}} \rangle \right)^2. \quad (7)$$

We shall consider the behavior of gradient flow-trained networks over the population loss induced by the limit of infinite training tasks/prompts $B \rightarrow \infty$:

$$L(\theta) = \lim_{B \rightarrow \infty} \hat{L}(\theta) = \frac{1}{2} \mathbb{E}_{w_\tau, x_{\tau,1}, \dots, x_{\tau,N}, x_{\tau,\text{query}}} \left[(\hat{y}_{\tau,\text{query}} - \langle w_\tau, x_{\tau,\text{query}} \rangle)^2 \right] \quad (8)$$

1 Introduction

2 Preliminaries

3 Main results

Convergence

Assumption 3 (Initialization) Let $\sigma > 0$ be a parameter, and let $\Theta \in \mathbb{R}^{d \times d}$ be any matrix satisfying $\|\Theta \Theta^\top\|_F = 1$ and $\Theta \Lambda \neq 0_{d \times d}$. We assume

$$W^{PV}(0) = \sigma \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & 1 \end{pmatrix}, \quad W^{KQ}(0) = \sigma \begin{pmatrix} \Theta \Theta^\top & 0_d \\ 0_d^\top & 0 \end{pmatrix}. \quad (10)$$

Theorem 4 (Convergence and limits) Consider gradient flow of a linear self-attention network f_{LSA} defined in (3) over the population loss (8). Suppose the initialization satisfies Assumption 3 with initialization scale $\sigma > 0$ satisfying $\sigma^2 \|\Gamma\|_{\text{op}} \sqrt{d} < 2$ where we have defined

$$\Gamma := \left(1 + \frac{1}{N}\right) \Lambda + \frac{1}{N} \text{tr}(\Lambda) I_d \in \mathbb{R}^{d \times d}.$$

Then gradient flow converges to a global minimum of the population loss (8). Moreover, W^{PV} and W^{KQ} converge to W_*^{PV} and W_*^{KQ} respectively, where

$$W_*^{KQ} = [\text{tr}(\Gamma^{-2})]^{-\frac{1}{4}} \begin{pmatrix} \Gamma^{-1} & 0_d \\ 0_d^\top & 0 \end{pmatrix}, \quad W_*^{PV} = [\text{tr}(\Gamma^{-2})]^{\frac{1}{4}} \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & 1 \end{pmatrix}. \quad (11)$$

Convergence

$$\begin{aligned}
 \hat{y}_{\text{query}} &= \begin{pmatrix} 0_d^\top & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{M} \sum_{i=1}^M x_i x_i^\top + \frac{1}{M} x_{\text{query}} x_{\text{query}}^\top & \frac{1}{M} \sum_{i=1}^M x_i x_i^\top w \\ \frac{1}{M} \sum_{i=1}^M w^\top x_i x_i^\top & \frac{1}{M} \sum_{i=1}^M w^\top x_i x_i^\top w \end{pmatrix} \begin{pmatrix} \Gamma^{-1} & 0_d \\ 0_d^\top & 0 \end{pmatrix} \begin{pmatrix} x_{\text{query}} \\ 0 \end{pmatrix} \\
 &= x_{\text{query}}^\top \Gamma^{-1} \left(\frac{1}{M} \sum_{i=1}^M x_i x_i^\top \right) w. \tag{12}
 \end{aligned}$$

- Obviously, when N is large enough, $\Gamma \rightarrow \Lambda$
- $\sum_{i=1}^M \frac{1}{M} x_i x_i^\top \rightarrow \Lambda$
- $\hat{y}_{\text{query}} \rightarrow x_{\text{query}}^\top w$

Convergence

$$\begin{aligned}
 \hat{y}_{\text{query}} &= \begin{pmatrix} 0_d^\top & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{M} \sum_{i=1}^M x_i x_i^\top + \frac{1}{M} x_{\text{query}} x_{\text{query}}^\top & \frac{1}{M} \sum_{i=1}^M x_i y_i \\ \frac{1}{M} \sum_{i=1}^M x_i^\top y_i & \frac{1}{M} \sum_{i=1}^M y_i^2 \end{pmatrix} \begin{pmatrix} \Gamma^{-1} & 0_d \\ 0_d^\top & 0 \end{pmatrix} \begin{pmatrix} x_{\text{query}} \\ 0 \end{pmatrix} \\
 &= x_{\text{query}}^\top \Gamma^{-1} \left(\frac{1}{M} \sum_{i=1}^M y_i x_i \right). \tag{13}
 \end{aligned}$$

- When h in test prompt is not linear, how's the thing going?
- Linear optimal as well!
- Suppose $(x_i, y_i) \sim D$ and $x_i \sim N(0, \Lambda)$
- $\Lambda^{-1} E_{(x, y) \sim D} [xy] = \arg \min_w E[(y - w^T x)^2]$
- Consider why?

Convergence

Theorem 5 Let \mathcal{D} be a distribution over $(x, y) \in \mathbb{R}^d \times \mathbb{R}$, whose marginal distribution on x is $\mathcal{D}_x = \mathcal{N}(0, \Lambda)$. Assume $\mathbb{E}_{\mathcal{D}}[y], \mathbb{E}_{\mathcal{D}}[xy], \mathbb{E}_{\mathcal{D}}[y^2 x x^\top]$ exist and are finite. Assume the test prompt is of the form $P = (x_1, y_1, \dots, x_M, y_M, x_{\text{query}})$, where $(x_i, y_i), (x_{\text{query}}, y_{\text{query}}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$. Let f_{LSA}^* be the LSA model with parameters W_*^{PV} and W_*^{KQ} in (11), and \hat{y}_{query} is the prediction for x_{query} given the prompt. If we define

$$a := \Lambda^{-1} \mathbb{E}_{(x,y) \sim \mathcal{D}} [xy], \quad \Sigma := \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[(xy - \mathbb{E}(xy)) (xy - \mathbb{E}(xy))^\top \right], \quad (15)$$

then, for $\Gamma = \Lambda + \frac{1}{N} \Lambda + \frac{1}{N} \text{tr}(\Lambda) I_d$. we have,

$$\begin{aligned} \mathbb{E}(\hat{y}_{\text{query}} - y_{\text{query}})^2 &= \underbrace{\min_{w \in \mathbb{R}^d} \mathbb{E}(\langle w, x_{\text{query}} \rangle - y_{\text{query}})^2}_{\text{Error of best linear predictor}} \\ &\quad + \boxed{\frac{1}{M}} \text{tr}[\Sigma \Gamma^{-2} \Lambda] + \boxed{\frac{1}{N^2}} \left[\|a\|_{\Gamma^{-2} \Lambda^3}^2 + 2 \text{tr}(\Lambda) \|a\|_{\Gamma^{-2} \Lambda^2}^2 + \text{tr}(\Lambda)^2 \|a\|_{\Gamma^{-2} \Lambda}^2 \right], \end{aligned} \quad (16)$$

where the expectation is over $(x_i, y_i), (x_{\text{query}}, y_{\text{query}}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$.

- The prediction error is at most $O(1/M + 1/N^2)$

Convergence

Let us now consider when \mathcal{D} corresponds to noiseless linear models, so that for some $w \in \mathbb{R}^d$, we have $(x, y) = (x, \langle w, x \rangle)$, in which case the prediction of the trained transformer is given by (12). Moreover, a simple calculation shows that the Σ from Theorem 5 takes the form $\Sigma = \|w\|_\Lambda^2 \Lambda + \Lambda w w^\top \Lambda$. Hence Theorem 5 implies the prediction error for the prompt $P = (x_1, \langle w, x_1 \rangle, \dots, x_M, \langle w, x_M \rangle, x_{\text{query}})$ is

$$\begin{aligned} & \mathbb{E}_{x_1, \dots, x_M, x_{\text{query}}} (\hat{y}_{\text{query}} - \langle w, x_{\text{query}} \rangle)^2 \\ &= \frac{1}{M} \left\{ \|w\|_{\Gamma^{-2}\Lambda^3}^2 + \text{tr}(\Gamma^{-2}\Lambda^2) \|w\|_\Lambda^2 \right\} \\ & \quad + \frac{1}{N^2} \left\{ \|w\|_{\Gamma^{-2}\Lambda^3}^2 + 2 \|w\|_{\Gamma^{-2}\Lambda^2}^2 \text{tr}(\Lambda) + \|w\|_{\Gamma^{-2}\Lambda}^2 \text{tr}(\Lambda)^2 \right\} \\ &\leq \frac{d+1}{M} \|w\|_\Lambda^2 + \frac{1}{N^2} \left[\|w\|_\Lambda^2 + 2 \|w\|_2^2 \text{tr}(\Lambda) + \|w\|_{\Lambda^{-1}}^2 \text{tr}(\Lambda)^2 \right]. \end{aligned}$$

The inequality above uses that $\Gamma \succ \Lambda$. Finally, if we assume that $w \sim \mathcal{N}(0, I_d)$ and denote κ as the condition number of Λ , then by taking expectations over w we get the following:

$$\begin{aligned} & \mathbb{E}_{x_1, \dots, x_M, x_{\text{query}}, w} (\hat{y}_{\text{query}} - \langle w, x_{\text{query}} \rangle)^2 \\ &\leq \frac{(d+1) \text{tr}(\Lambda)}{M} + \frac{1}{N^2} [\text{tr}(\Lambda) + 2d \text{tr}(\Lambda) + \text{tr}(\Lambda^{-1}) \text{tr}(\Lambda)^2] \\ &\leq \frac{(d+1) \text{tr}(\Lambda)}{M} + \frac{(1 + 2d + d^2 \kappa) \text{tr}(\Lambda)}{N^2}, \end{aligned}$$

- $\text{tr}(\Lambda)$, condition number ($\kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$), covariate dimension d can also affect the precision of prediction

Prompts with Random Covariate Distributions

Covariate Shifts

When $D_x^{train} \neq D_x^{test}$, the approximation does not hold when M and N are large. For instance, if we consider test prompts where the covariates are scaled by a constant $c \neq 1$:

$$\hat{y}_{query} \rightarrow x_{query}^T \Lambda^{-1} c^2 \Lambda w \neq x_{query}^T w$$

- To further research on this problem, we set the covariate matrix Λ to random
- Then, for each task τ and coordinate $i \in [d]$, we sample $\lambda_{\tau,i}$ independently such that the distribution of each $\lambda_{\tau,i}$ is fixed and has finite third moments and is strictly positive almost surely. We then form a diagonal matrix:

$$\Lambda_\tau = \text{diag}(\lambda_{\tau,1}, \dots, \lambda_{\tau,d})$$

Prompts with Random Covariate Distributions

Theorem 8 (Global convergence with random covariance) *Consider gradient flow of the linear self-attention network f_{LSA} defined in (3) over the population loss (20), where Λ_τ are diagonal with independent diagonal entries which are strictly positive a.s. and have finite third moments. Suppose the initialization satisfies Assumption 3, $\|\mathbb{E}\Lambda_\tau\Theta\|_F \neq 0$, with initialization scale $\sigma > 0$ satisfying*

$$\sigma^2 < \frac{2 \|\mathbb{E}\Lambda_\tau\Theta\|_F^2}{\sqrt{d} \left[\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 \right]}. \quad (21)$$

Then gradient flow converges to a global minimum of the population loss (20). Moreover, W^{PV} and W^{KQ} converge to W_^{PV} and W_*^{KQ} respectively, where*

$$\begin{aligned} W_*^{KQ} &= \left\| [\mathbb{E}\Gamma_\tau\Lambda_\tau^2]^{-1} \mathbb{E} [\Lambda_\tau^2] \right\|_F^{-\frac{1}{2}} \cdot \begin{pmatrix} [\mathbb{E}\Gamma_\tau\Lambda_\tau^2]^{-1} [\mathbb{E}\Lambda_\tau^2] & 0_d \\ 0_d^\top & 0 \end{pmatrix}, \\ W_*^{PV} &= \left\| [\mathbb{E}\Gamma_\tau\Lambda_\tau^2]^{-1} \mathbb{E} [\Lambda_\tau^2] \right\|_F^{\frac{1}{2}} \cdot \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & 1 \end{pmatrix}, \end{aligned} \quad (22)$$

where $\Gamma_\tau = \frac{N+1}{N}\Lambda_\tau + \frac{1}{N} \text{tr}(\Lambda_\tau)I_d \in \mathbb{R}^{d \times d}$ and the expectations above are over the distribution of Λ_τ .

Prompts with Random Covariate Distributions

$$\begin{aligned}
 & \hat{y}_{\text{query}} \\
 &= \begin{pmatrix} 0_d^\top & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{M} \sum_{i=1}^M x_i x_i^\top + \frac{1}{M} x_{\text{query}} x_{\text{query}}^\top & \frac{1}{M} \sum_{i=1}^M x_i y_i \\ \frac{1}{M} \sum_{i=1}^M x_i^\top y_i & \frac{1}{M} \sum_{i=1}^M y_i^2 \end{pmatrix} \begin{pmatrix} [\mathbb{E} \Gamma_\tau \Lambda_\tau^2]^{-1} \mathbb{E} \Lambda_\tau^2 & 0_d \\ 0_d^\top & 0 \end{pmatrix} \begin{pmatrix} x_{\text{query}} \\ 0 \end{pmatrix} \\
 &= x_{\text{query}}^\top \cdot [\mathbb{E} \Lambda_\tau^2] [\mathbb{E} \Gamma_\tau \Lambda_\tau^2]^{-1} \cdot \left[\frac{1}{M} \sum_{i=1}^M x_i x_i^\top \right] w \\
 &\rightarrow x_{\text{query}}^\top \cdot [\mathbb{E} \Lambda_\tau^2] [\mathbb{E} \Gamma_\tau \Lambda_\tau^2]^{-1} \cdot \Lambda_{\text{new}} w \quad \text{almost surely when } M \rightarrow \infty. \tag{23}
 \end{aligned}$$

- One clear example is: considering $\lambda_{\tau,i} \sim \text{Exp}(1)$ and $\Gamma_\tau \rightarrow \Lambda_\tau$, then $E(\Lambda_\tau) = 1$, $E(\Lambda_\tau^2) = 2$, $E(\Lambda_\tau^3) = 6$, as a result:

$$E[\hat{y}_{\text{query}} | x_{\text{query}}, w] \rightarrow \frac{1}{3} w^\top x_{\text{query}}$$

Experiments with Large, Nonlinear Transformers

