
Transformers as Statisticians: Provable In-Context Learning with In-Context Algorithm Selection

Yu Bai*

Salesforce Research
yu.bai@salesforce.com

Fan Chen*

Massachusetts Institute of Technology
fanchen@mit.edu

Huan Wang

Salesforce Research
huan.wang@salesforce.com

Caiming Xiong

Salesforce Research
cxiong@salesforce.com

Song Mei*

UC Berkeley
songmei@berkeley.edu

Abstract

Neural sequence models based on the transformer architecture have demonstrated remarkable *in-context learning* (ICL) abilities, where they can perform new tasks when prompted with training and test examples, without any parameter update to the model. This work first provides a comprehensive statistical theory for transformers to perform ICL. Concretely, we show that transformers can implement a broad class of standard machine learning algorithms in context, such as least squares, ridge regression, Lasso, learning generalized linear models, and gradient descent on two-layer neural networks, with near-optimal predictive power on various in-context data distributions. Using an efficient implementation of in-context gradient descent as the underlying mechanism, our transformer constructions admit mild size bounds, and can be learned with polynomially many pretraining sequences.

Building on these “base” ICL algorithms, intriguingly, we show that transformers can implement more complex ICL procedures involving *in-context algorithm selection*, akin to what a statistician can do in real life—A *single* transformer can adaptively select different base ICL algorithms—or even perform qualitatively different tasks—on different input sequences, without any explicit prompting of the right algorithm or task. We both establish this in theory by explicit constructions, and also observe this phenomenon experimentally. In theory, we construct two general mechanisms for algorithm selection with concrete examples: pre-ICL testing, and post-ICL validation. As an example, we use the post-ICL validation mechanism to construct a transformer that can perform nearly Bayes-optimal ICL on a challenging task—noisy linear models with mixed noise levels. Experimentally, we demonstrate the strong in-context algorithm selection capabilities of standard transformer architectures.

1 Introduction

Large neural sequence models have demonstrated remarkable *in-context learning* (ICL) capabilities [12], where models can make accurate predictions on new tasks when prompted with training examples from the same task, in a zero-shot fashion without any parameter update to the model. A prevalent example is large language models based on the transformer architecture [84], which can perform a diverse range of tasks in context when trained on enormous text [12, 90]. Recent models

*Equal technical and directional contributions.

Code is available at <https://github.com/allenbai01/transformers-as-statisticians>.

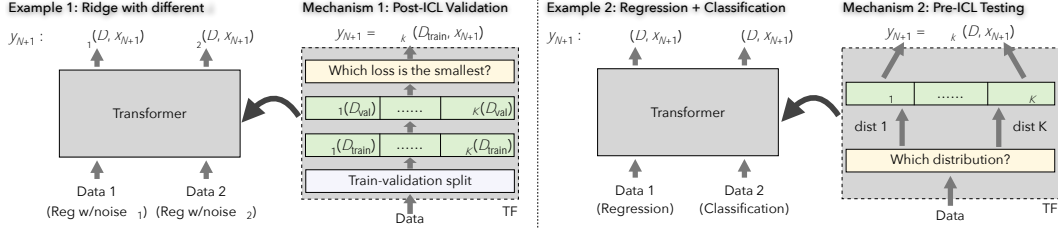


Figure 1: **Illustration of in-context algorithm selection, and two mechanisms constructed in our theory.** *Left, middle-left:* A single transformer can perform ridge regression with different λ 's on input sequences with different observation noise; we prove this by the **post-ICL validation** mechanism (Section 4.1). *Middle-right, right:* A single transformer can perform linear regression on regression data and logistic regression on classification data; we prove this via the **pre-ICL testing** mechanism (Section 4.2).

in this paradigm such as GPT-4 achieve surprisingly impressive ICL performance that makes them akin to a general-purpose agent in many aspects [65, 14]. Such strong capabilities call for better understandings, which a recent line of work tackles from various aspects [49, 94, 28, 72, 15, 57, 64].

Recent pioneering work of Garg et al. [31] proposes an interpretable and theoretically amenable setting for understanding ICL in transformers. They perform ICL experiments where input tokens are real-valued (input, label) pairs generated from standard statistical models such as linear models (and the sparse version), neural networks, and decision trees. Garg et al. [31] find that transformers can learn to perform ICL with prediction power (and fitted functions) matching standard machine learning algorithms for these settings, such as least squares for linear models, and Lasso for sparse linear models. Subsequent work further studies the internal mechanisms [2, 86, 18], expressive power [2, 32], and generalization [47] of transformers in this setting. However, these works only showcase simple mechanisms such as regularized regression [31, 2, 47] or gradient descent [2, 86, 18], which are arguably only a small subset of what transformers are capable of in practice; or expressing universal function classes not specific to ICL [89, 32]. This motivates the following question:

How do transformers learn in context beyond implementing simple algorithms?

This paper makes steps on this question by making two main contributions: (1) We **unveil a general mechanism—in-context algorithm selection**—by which a *single* transformer can adaptively *select different “base” ICL algorithms* to use on *different ICL instances*, without any explicit prompting of the right algorithm to use in the input sequence. For example, a transformer may choose to perform ridge regression with regularization λ_1 on ICL instance 1, and λ_2 on ICL instance 2 (Figure 2); or perform regression on ICL instance 1 and classification on ICL instance 2 (Figure 5). This adaptivity allows transformers to achieve much stronger ICL performance than the base ICL algorithms. We both prove this in theory, and demonstrate this phenomenon empirically on standard transformer architectures. (2) Along the way, equally importantly, we present a comprehensive theory for ICL in transformers by establishing end-to-end quantitative guarantees for the **expressive power, in-context prediction performance, and sample complexity of pretraining**. These results add upon the recent line of work on the statistical learning theory of transformers [97, 89, 27, 39], and lay out a foundation for the intriguing special case where the *learning targets are themselves ICL algorithms*.

A detailed summary of our contributions is as follows.

- We prove that transformers can implement a broad class of standard machine learning algorithms in context, such as least squares, ridge regression, Lasso, convex risk minimization for learning generalized linear models (such as logistic regression), and gradient descent for two-layer neural networks (Section 3). Our constructions admit mild bounds on the number of layers, heads, and weight norms, and achieve near-optimal prediction power on many in-context data distributions.
- Technically, the above transformer constructions build on a new efficient implementation of in-context gradient descent (Appendix D), which could be broadly applicable. For a broad class of smooth convex empirical risks over the in-context training data, we construct an $(L + 1)$ -layer transformer that approximates L steps of gradient descent. Notably, the approximation error accumulates only *linearly* in L , utilizing a stability-like property of smooth convex optimization.
- We prove that transformers can perform in-context algorithm selection (Section 4). We construct two algorithm selection mechanisms: Post-ICL validation (Section 4.1), and Pre-ICL testing

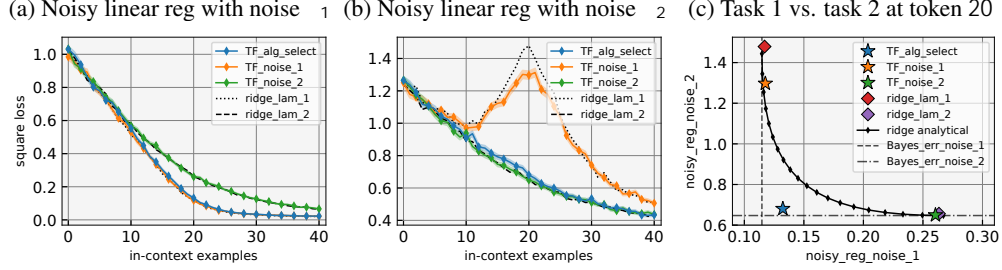


Figure 2: In-context algorithm selection on two separate noisy linear regression tasks with noise $(\sigma_1; \sigma_2) = (0.1; 0.5)$. (a,b) A **single transformer** `TF_alg_select` **simultaneously approaches the performance of the two individual Bayes predictors** `ridge_lam_1` on task 1 and `ridge_lam_2` on task 2. (c) At token 20 (using example `f0;:::;19g` for training), `TF_alg_select` approaches the Bayes error on two tasks simultaneously, and **outperforms ridge regression with any fixed** λ . (a,b,c) Note that transformers pretrained on a single task (`TF_noise_1`, `TF_noise_2`) perform near-optimally on that task but suboptimally on the other task. More details about the setup and training method can be found in Appendix M.2.

(Section 4.2). For both mechanisms, we provide general constructions as well as concrete examples. Figure 1 provides a pictorial illustration of the two mechanisms.

- As a concrete application, using the post-ICL validation mechanism, we construct a transformer that can perform nearly Bayes-optimal ICL on noisy linear models with *mixed* noise levels (Section 4.1.1), a more complex task than those considered in existing work.
- We provide the first line of results for *pretraining* transformers to perform the various ICL tasks above, from polynomially many training sequences (Section 5 & Appendix K).
- Experimentally, we find that learned transformers indeed exhibit strong in-context algorithm selection capabilities in the settings considered in our theory (Section 6). For example, Figure 2 shows that a *single* transformer can approach the individual Bayes risks (the optimal risk among all possible algorithms) simultaneously on two noisy linear models with different noise levels.

Transformers as statisticians We humbly remark that the typical toolkit of a statistician contains much more beyond those covered in this work, including and not limited to inference, uncertainty quantification, and theoretical analysis. This work merely aims to show the algorithm selection capability of transformers, akin to what a statistician *can* do.

Related work Our work is intimately related to the lines of work on in-context learning, theoretical understandings of transformers, as well as other formulations for learning-to-learn such as meta-learning. Due to limited space, we discuss these related work in Appendix A.

2 Preliminaries

We consider a sequence of N input vectors $\{h_i\}_{i=1}^N \in \mathbb{R}^D$, written compactly as an input matrix $\mathbf{H} = [h_1; \dots; h_N] \in \mathbb{R}^{D \times N}$, where each h_i is a column of \mathbf{H} (also a *token*). Throughout this paper, we let $\text{ReLU}(t) := \max\{t, 0\}$ denote the standard relu activation.

2.1 Transformers

We consider transformer architectures that process any input sequence $\mathbf{H} \in \mathbb{R}^{D \times N}$ by applying (encoder-mode²) attention layers and MLP layers formally defined as follows.

Definition 1 (Attention layer). A (self-)attention layer with M heads is denoted as $\text{Attn}(\cdot)$ with parameters $\{V_m, Q_m, K_m\}_{m \in [M]} \in \mathbb{R}^{D \times D}$. On any input sequence $\mathbf{H} \in \mathbb{R}^{D \times N}$,

$$\hat{\mathbf{H}} = \text{Attn}(\mathbf{H}) := \mathbf{H} + \frac{1}{N} \prod_{m=1}^M (V_m \mathbf{H}) \quad (Q_m \mathbf{H})^\top (K_m \mathbf{H}) \in \mathbb{R}^{D \times N}; \quad (1)$$

where $\text{ReLU} : \mathbb{R} \rightarrow \mathbb{R}$ is the ReLU function. In vector form,

$$\hat{h}_i = [\text{Attn}(\mathbf{H})]_i = h_i + \prod_{m=1}^M \frac{1}{N} \prod_{j=1}^N (h_j Q_m h_i; K_m h_j i) V_m h_j;$$

²Many of our results can be generalized to decoder-based architectures; see Appendix C for a discussion.

Above, (1) uses a normalized ReLU³ activation $t \mathbb{1}(t) = N$ in place of the standard softmax activation; we remark this activation is also found to work well empirically in recent studies [78, 93].

Definition 2 (MLP layer). A (token-wise) MLP layer with hidden dimension D' is denoted as $\text{MLP}(\cdot)$ with parameters $(\mathbf{W}_1; \mathbf{W}_2) \in \mathbb{R}^{D' \times D} \times \mathbb{R}^{D \times D'}$. On any input sequence $\mathbf{H} \in \mathbb{R}^{D \times N}$,

$$\hat{\mathbf{H}} = \text{MLP}(\mathbf{H}) := \mathbf{H} + \mathbf{W}_2 (\mathbf{W}_1 \mathbf{H});$$

where $\mathbb{1} : \mathbb{R} \rightarrow \mathbb{R}$ is the ReLU function. In vector form, we have $\hat{\mathbf{h}}_i = \mathbf{h}_i + \mathbf{W}_2 (\mathbf{W}_1 \mathbf{h}_i)$.

We consider a transformer architecture with $L \geq 1$ transformer layers, each consisting of a self-attention layer followed by an MLP layer.

Definition 3 (Transformer). An L -layer transformer, denoted as $\text{TF}(\cdot)$, is a composition of L self-attention layers each followed by an MLP layer: $\mathbf{H}^{(L)} = \text{TF}(\mathbf{H}^{(0)})$, where $\mathbf{H}^{(0)} \in \mathbb{R}^{D \times N}$ is the input sequence, and

$$\mathbf{H}^{(\ell)} = \text{MLP}_{\text{mlp}}^{(\ell)} \text{Attn}_{\text{attn}}^{(\ell)} \mathbf{H}^{(\ell-1)}; \quad \ell \in \{1, \dots, L\};$$

Above, the parameter $\text{Attn}_{\text{attn}}^{(\ell)} = (Q_{\text{attn}}^{(\ell)}; K_{\text{attn}}^{(\ell)}; V_{\text{attn}}^{(\ell)})$ consists of the attention layers $f_{\text{attn}}^{(\ell)} = f(\mathbf{Q}_m^{(\ell)}; \mathbf{K}_m^{(\ell)}; \mathbf{V}_m^{(\ell)})_{m \in [M]} \in \mathbb{R}^{D \times D}$ and the MLP layers $\text{MLP}_{\text{mlp}}^{(\ell)} = (\mathbf{W}_1^{(\ell)}; \mathbf{W}_2^{(\ell)}) \in \mathbb{R}^{D' \times D} \times \mathbb{R}^{D \times D'}$. We will frequently consider ‘‘attention-only’’ transformers with $\mathbf{W}_1^{(\ell)}; \mathbf{W}_2^{(\ell)} = \mathbf{0}$, which we denote as $\text{TF}^0(\cdot)$ for shorthand, with $\text{Attn}_{\text{attn}}^{(\ell)} = (Q_{\text{attn}}^{(\ell)}; K_{\text{attn}}^{(\ell)}; V_{\text{attn}}^{(\ell)})$.

We additionally define the following norm of a transformer $\text{TF}(\cdot)$:

$$\|\text{TF}(\cdot)\| := \max_{\ell \in [L]} \left(\max_{m \in [M]} \left(k_{\text{Q}_m^{(\ell)}} + k_{\text{K}_m^{(\ell)}} \right) + \sum_{m=1}^M \left(k_{\text{V}_m^{(\ell)}} + k_{\text{W}_1^{(\ell)}} + k_{\text{W}_2^{(\ell)}} \right) \right); \quad (2)$$

In (2), the choices of the operator norm and max/sums are for convenience only and not essential, as our results (e.g. for pretraining) depend only logarithmically on $\|\text{TF}(\cdot)\|$.

2.2 In-context learning

In an in-context learning (ICL) instance, the model is given a dataset $D = \{(\mathbf{x}_i; y_i)_{i \in [N]}\}$ and a new test input $\mathbf{x}_{N+1} \in \mathbb{R}^d$ for some data distribution \mathcal{P} , where $\{\mathbf{x}_i; y_i\}_{i \in [N]} \in \mathbb{R}^d$ are the input vectors, $\{y_i\}_{i \in [N]} \in \mathbb{R}$ are the corresponding labels (e.g. real-valued for regression, or $\{0, 1\}$ -valued for binary classification), and \mathbf{x}_{N+1} is the test input on which the model is required to make a prediction. Different from standard supervised learning, in ICL, each instance $(D; \mathbf{x}_{N+1})$ is in general drawn from a different distribution \mathcal{P}_j , such as a linear model with a new ground truth coefficient $\mathbf{w}_{?j} \in \mathbb{R}^d$. Our goal is to construct *fixed* transformer to perform ICL on a large set of \mathcal{P}_j 's.

We consider using transformers to perform ICL, in which we encode $(D; \mathbf{x}_{N+1})$ into an input sequence $\mathbf{H} \in \mathbb{R}^{D \times (N+1)}$. In our theory, we use the following format, where the first two rows contain $(D; \mathbf{x}_{N+1})$ (zero at the location for y_{N+1}), and the third row contains fixed vectors $\{\mathbf{p}_i\}_{i \in [N+1]}$ with ones, zeros, and indicator for being the train token (similar to a positional encoding vector):

$$\mathbf{H} = \begin{matrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_N & \mathbf{x}_{N+1} \\ y_1 & y_2 & \dots & y_N & 0 \\ \mathbf{p}_1 & \mathbf{p}_2 & \dots & \mathbf{p}_N & \mathbf{p}_{N+1} \end{matrix} \in \mathbb{R}^{D \times (N+1)}; \quad \mathbf{p}_i := \begin{matrix} 0_{D-(d+3)} \\ 1 \\ 1_{fi < N+1g} \end{matrix} \in \mathbb{R}^{D-(d+1)}; \quad (3)$$

We will choose $D = \Theta(d)$, so that the hidden dimension of \mathbf{H} is at most a constant multiple of d . We then feed \mathbf{H} into a transformer to obtain the output $\hat{\mathbf{H}} = \text{TF}(\mathbf{H}) \in \mathbb{R}^{D \times (N+1)}$ with the same shape, and *read out* the prediction \hat{y}_{N+1} from the $(d+1; N+1)$ -th entry of $\hat{\mathbf{H}} = [\hat{\mathbf{h}}_i]_{i \in [N+1]}$ (the entry corresponding to the missing test label): $\hat{y}_{N+1} = \text{read}_y(\hat{\mathbf{H}}) := (\hat{\mathbf{h}}_{N+1})_{d+1}$. The goal is to predict

³For each query index i , the attention weights $f(\mathbf{h}_i; \mathbf{Q}_m; \mathbf{h}_j) = \sum_{j \in [N]} g_j$ is also a set of non-negative weights that sum to $O(1)$ (similar as a softmax probability distribution) in typical scenarios. Also, our approximation results can potentially be generalized to softmax attention e.g. using the technique of [32].

\hat{y}_{N+1} that is close to y_{N+1} $\mathbb{P}_{y|\mathbf{x}_{N+1}}$ measured by proper losses. We emphasize that we consider predicting only at the last token \mathbf{x}_{N+1} , which is without much loss of generality.⁴

Miscellaneous setups We assume bounded features and labels throughout the paper (unless otherwise specified, e.g. when \mathbf{x}_i is Gaussian): $k\mathbf{x}_i k_2 \leq B_x$ and $|y_i| \leq B_y$ with probability one. We use the standard notation $\mathbf{X} = [\mathbf{x}_1^\top; \dots; \mathbf{x}_N^\top] \in \mathbb{R}^{N \times d}$ and $\mathbf{y} = [y_1; \dots; y_N] \in \mathbb{R}^N$ to denote the matrix of inputs and vector of labels, respectively. To prevent the transformer from blowing up on tail events, in all our results concerning (statistical) in-context prediction powers, we consider a clipped prediction $\hat{y}_{N+1} = \text{read}_y(\hat{\mathbf{H}}) := \text{clip}_R(\hat{\mathbf{H}}_{N+1})_{d+1}$, where $\text{clip}_R(t) := \text{Proj}_{[-R; R]}(t)$ is the standard clipping operator with (a suitably large) radius $R \geq 0$ that varies in different problems.

3 Basic in-context learning algorithms

We begin by constructing transformers that approximately implement a variety of standard machine learning algorithms in context, with mild size bounds and near-optimal prediction power on many standard in-context data distributions.

3.1 In-context ridge regression and least squares

Consider the standard ridge regression estimator over the in-context training examples D with regularization $\lambda > 0$ (reducing to least squares at $\lambda = 0$ and $N \rightarrow \infty$):

$$\mathbf{w}_{\text{ridge}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2N} \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2; \quad (\text{ICRidge})$$

We show that transformers can approximately implement (ICRidge) (proof in Appendix F.1).

Theorem 4 (Implementing in-context ridge regression). *For any $\epsilon > 0, \lambda > 0$ with $\lambda := \frac{\epsilon}{B_x + B_y}$, $B_w > 0$, and $\lambda < B_x B_w = 2$, there exists an L -layer attention-only transformer TF^0 with*

$$L = d \log(B_x B_w = 2) \epsilon^{-1}; \quad \max_{i \in [L]} M^{(i)} \leq 3; \quad \forall i \in [L] \quad 4R + 8(\epsilon + \lambda)^{-1}; \quad (4)$$

(with $R := \max\{B_x B_w; B_y; 1\}g$) such that the following holds. On any input data $(D; \mathbf{x}_{N+1})$ such that the problem (ICRidge) is well-conditioned and has a bounded solution:

$$\min(\mathbf{X}^\top \mathbf{X} = N) \leq \max(\mathbf{X}^\top \mathbf{X} = N) \leq \frac{1}{\lambda}; \quad \mathbf{w}_{\text{ridge}} \in \mathbb{R}^d; \quad B_w = 2; \quad (5)$$

TF^0 approximately implements (ICRidge): The prediction $\hat{y}_{N+1} = \text{read}_y(\text{TF}^0(\mathbf{H}))$ satisfies

$$|\hat{y}_{N+1} - \mathbf{w}_{\text{ridge}}^\top \mathbf{x}_{N+1}| \leq \epsilon; \quad (6)$$

Theorem 4 presents the first quantitative construction for end-to-end in-context ridge regression up to arbitrary precision, and improves upon Akyürek et al. [2] whose construction does not give (or directly imply) an explicit error bound like (6). Further, the bounds on the number of layers and heads in (4) are mild (constant heads and logarithmically many layers).

Near-optimal in-context prediction power for linear problems Combining Theorem 4 with standard analyses of linear regression yields the following corollaries (proofs in Appendix F.3 & F.4).

Corollary 5 (Near-optimal linear regression with transformers by approximating least squares). *For any $N \geq \Theta(d)$, there exists an $O(\log(N =))$ -layer transformer \mathcal{T} , such that on any \mathbb{P} satisfying standard statistical assumptions for least squares (Assumption A), its ICL prediction \hat{y}_{N+1} achieves*

$$\mathbb{E}_{(D; \mathbf{x}_{N+1}; y_{N+1}) \sim \mathbb{P}} [(\hat{y}_{N+1} - y_{N+1})^2] \leq \inf_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}; y) \sim \mathbb{P}} (y - \mathbf{w}^\top \mathbf{x})^2 + \Theta(d^{-2} = N);$$

Assumption A requires only generic tail properties such as sub-Gaussianity, and *not* realizability (i.e., \mathbb{P} follows a true linear model); κ ; β above denote the covariance condition number and the

⁴Our constructions may be generalized to predicting at every token, by using a decoder architecture and potentially different input formats correspondingly (cf. Appendix C). Our theory focuses on predicting at the last token only, which simplifies the setting. Our experiments test both settings.

noise level therein. The $\Theta(d^{-2}=N)$ excess risk is known to be rate-optimal for linear regression [38], and Corollary 5 achieves this in context with a transformer with only logarithmically many layers.

Next, consider Bayesian linear models where each in-context data distribution $P = P_{\mathbf{w}_\gamma}^{\text{lin}}$ is drawn from a Gaussian prior $\mathbf{w}_\gamma \sim N(0; I_{d=1})$, and $(\mathbf{x}; y) \sim P_{\mathbf{w}_\gamma}^{\text{lin}}$ is sampled as $\mathbf{x} \sim N(0; I_d)$, $y = \langle \mathbf{w}_\gamma; \mathbf{x} \rangle + N(0; \sigma^2)$. It is a standard result that the Bayes estimator of y_{N+1} given $(D; \mathbf{x}_{N+1})$ is given by ridge regression (ICRidge): $\hat{y}_{N+1}^{\text{Bayes}} := \langle \mathbf{w}_{\text{ridge}}; \mathbf{x}_{N+1} \rangle$ with $\sigma^2 = d^{-2}=N$. We show that transformers achieve nearly-Bayes risk for this problem, and we use

$$\text{BayesRisk} := E_{\mathbf{w}_\gamma \sim P; (D; \mathbf{x}_{N+1}; y_{N+1}) \sim P_{\mathbf{w}_\gamma}^{\text{lin}}} \frac{1}{2} (\hat{y}_{N+1}^{\text{Bayes}} - y_{N+1})^2$$

to denote the Bayes risk of this problem under prior P .

Corollary 6 (Nearly-Bayes linear regression with transformers by approximating ridge regression). *Under the Bayesian linear model above with $N \geq \max\{d=10; O(\log(1/\sigma^2))\}$, there exists a $L = O(\log(1/\sigma^2))$ -layer transformer such that $E_{\mathbf{w}_\gamma; (D; \mathbf{x}_{N+1}; y_{N+1})} \frac{1}{2} (\hat{y}_{N+1} - y_{N+1})^2 \leq \text{BayesRisk} + \sigma^2$.*

Generalized linear models In Appendix G, we extend the above results to generalized linear models [53] and show that transformers can approximate the corresponding convex risk minimization algorithm in context (which includes logistic regression for linear classification as an important special case), and achieve near-optimal excess risk under standard statistical assumptions.

3.2 In-context Lasso

Consider the standard Lasso estimator [82] which minimizes an ℓ_1 -regularized linear regression loss $\hat{\mathcal{L}}_{\text{lasso}}$ over the in-context training examples D :

$$\mathbf{w}_{\text{lasso}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} \hat{\mathcal{L}}_{\text{lasso}}(\mathbf{w}) = \frac{1}{2N} \sum_{i=1}^N (\langle \mathbf{w}; \mathbf{x}_i \rangle - y_i)^2 + \lambda \|\mathbf{w}\|_1; \quad (\text{ICLasso})$$

We show that transformers can also approximate in-context Lasso with a mild number of layers, and can perform sparse linear regression in standard sparse linear models (proofs in Appendix H).

Theorem 7 (Implementing in-context Lasso). *For any $N \geq 0$, $\lambda > 0$, $B_w > 0$, and $\sigma > 0$, there exists a L -layer transformer TF with*

$$L = \lceil B_w^2 \sigma^{-2} \rceil + 1; \quad \max_{\mathbf{x} \in [L]} M(\mathbf{x}) \leq 2; \quad \max_{\mathbf{x} \in [L]} D(\mathbf{x}) \leq 2d; \quad \|\mathbf{J}\| \|\mathbf{J}\| \leq O(R + (1 + N)^{-1})$$

(where $R := \max\{B_x B_w; B_y; 1\}$) such that the following holds. On any input data $(D; \mathbf{x}_{N+1})$ such that $\max(\mathbf{X}^\top \mathbf{X} = N)$ and $\|\mathbf{w}_{\text{lasso}}\|_2 \leq B_w = 2$, TF (H⁽⁰⁾) approximately implements (ICLasso), in that it outputs $\hat{y}_{N+1} = \langle \mathbf{x}_{N+1}; \hat{\mathbf{w}} \rangle$ with $\hat{\mathcal{L}}_{\text{lasso}}(\hat{\mathbf{w}}) \leq \hat{\mathcal{L}}_{\text{lasso}}(\mathbf{w}_{\text{lasso}}) + \sigma^2$.

Theorem 8 (Near-optimal sparse linear regression with transformers by approximating Lasso). *For any $d; N \geq 1$; $\lambda > 0$; $B_w^2 \sigma^{-2} > 0$, there exists a $\Theta((B_w^2 \sigma^{-2})^2 (1 + (d=N)))$ -layer transformer such that the following holds: For any s and $N \geq O(s \log(d/\sigma))$, suppose that P is an s -sparse linear model: $\mathbf{x}_i \sim N(0; I_d)$, $y_i = \langle \mathbf{w}_\gamma; \mathbf{x}_i \rangle + N(0; \sigma^2)$ for any $\|\mathbf{w}_\gamma\|_2 \leq B_w = 2$ and $\|\mathbf{w}_\gamma\|_0 \leq s$, then with probability at least $1 - \epsilon$ (over the randomness of D), the transformer output \hat{y}_{N+1} achieves*

$$E_{(\mathbf{x}_{N+1}; y_{N+1}) \sim P} (\hat{y}_{N+1} - y_{N+1})^2 \leq \sigma^2 [1 + O(s \log(d/\sigma) = N)];$$

The $\Theta(s \log d = N)$ excess risk obtained in Theorem 8 is optimal up to log factors [62, 87]. We remark that Theorem 8 is not a direct corollary of Theorem 7; Rather, the bound on the number of layers in Theorem 8 requires a sharper convergence analysis of the (ICLasso) problem under sparse linear models (Appendix H.2), similar to [1].

3.3 Proof technique: In-context gradient descent

The constructions in Section 3.1 and 3.2 is built on the following result for approximating in-context (proximal) gradient descent on (regularized) convex losses.

Theorem 9 (ICGD; Informal version of Theorem D.1 & D.2). *For a broad class of convex losses of form $\hat{\mathcal{L}} = \frac{1}{N} \sum_{i=1}^N (\langle \mathbf{w}; \mathbf{x}_i \rangle - y_i)^2 + R(\mathbf{w})$, there exists an L -layer transformer that takes in any $(D; \mathbf{w}^0)$ and outputs $\hat{\mathbf{w}}^L$ such that $\|\hat{\mathbf{w}}^L - \mathbf{w}_{\{\text{GD}; \text{PGD}\}}^L\|_2 \leq O(L^{-\alpha})$, by composing L identical layers each $O(\alpha)$ -approximating a single step of GD (so that $O(L^{-\alpha})$ is a linear error accumulation).*

Theorem 9 is established in two main steps:

- Approximating one-step of ICGD using one attention layer (Proposition E.1), which substantially generalizes that of von Oswald et al. [86] (which only does GD on square losses with a *linear* self-attention), and is simpler than the ones in Akyürek et al. [2] and Giannou et al. [32].
- Stacking L of the above layer to approximate L steps of ICGD. Done naively, the error accumulation of this stacking operation is exponential in L in the worst case. We utilize the stability of *convex* gradient descent (Lemma D.1) to obtain the *linear* in L error accumulation in Theorem 9.

In Appendix D.3, we also give results for *non-convex* GD on two-layer neural nets, though with a worse (exponential in L) error accumulation as expected.

4 In-context algorithm selection

We now show that transformers can perform various kinds of *in-context algorithm selection*, which allows them to implement more complex ICL procedures by adaptively selecting different “base” algorithms on different input sequences. We construct two general mechanisms: *Post-ICL validation*, and *Pre-ICL testing*; See Figure 1 for a pictorial illustration.

4.1 Post-ICL validation mechanism

In our first mechanism, post-ICL validation, the transformer begins by implementing a *train-validation split* $D = (D_{\text{train}}; D_{\text{val}})$, and running K base ICL algorithms on D_{train} . Let $f_k g_{k \in [K]} : (\mathbb{R}^d \rightarrow \mathbb{R})$ denote the K learned predictors, and

$$\mathbb{E}_{\text{val}}(f) := \frac{1}{|D_{\text{val}}|} \mathbb{P}_{(x_i; y_i) \in D_{\text{val}}} (f(x_i); y_i) \quad (7)$$

denote the validation loss of any predictor f .

We show that (proof in Appendix I.1) a 3-layer transformer can output a predictor \hat{f} that achieves nearly the smallest validation loss, and thus nearly optimal expected loss if \mathbb{E}_{val} concentrates around the expected loss L . Below, the input sequence \mathbf{H} uses a generalized positional encoding $\mathbf{p}_i := [\mathbf{0}_{D-(d+3)}; 1; t_i]$ in (3), where $t_i := 1$ for $i \in D_{\text{train}}$, $t_i := -1$ for $i \in D_{\text{val}}$, and $t_{N+1} := 0$.

Proposition 10 (In-context algorithm selection via train-validation split). *Suppose that $(;)$ in (7) is approximable by sum of relus (Definition D.1, which includes all C^3 -smooth bivariate functions). Then there exists a 3-layer transformer TF that maps (defining $y'_i = y_i \mathbb{1}_{f_i < N+1}$)*

$$\mathbf{h}_i = [\mathbf{x}_i; y'_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbf{0}_{K+1}; 1; t_i] \quad \mathbf{h}'_i = [\mathbf{x}_i; y'_i; \hat{f}(\mathbf{x}_i); 1; t_i]; \quad i \in [N+1];$$

where the predictor $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex combination of f_k : $\mathbb{E}_{\text{val}}(\hat{f}_k) = \min_{k \in [K]} \mathbb{E}_{\text{val}}(f_k) + g$. As a corollary, for any convex risk $L : (\mathbb{R}^d \rightarrow \mathbb{R}) \rightarrow \mathbb{R}$, \hat{f} satisfies

$$L(\hat{f}) = \min_{k \in [K]} L(f_k) + \max_{k \in [K]} \mathbb{E}_{\text{val}}(f_k) - L(f_k) + g$$

Ridge regression with in-context regularization selection As an example, we use Proposition 10 to construct a transformer to perform in-context ridge regression with regularization selection according to the *unregularized* validation loss $\mathbb{E}_{\text{val}}(\mathbf{w}) := \frac{1}{2|D_{\text{val}}|} \mathbb{P}_{(x_i; y_i) \in D_{\text{val}}} (\langle \mathbf{w}; \mathbf{x}_i \rangle - y_i)^2$ (proof in Appendix I.2). Let $\lambda_1; \dots; \lambda_K \geq 0$ be K fixed regularization strengths.

Theorem 11 (Ridge regression with in-context regularization selection). *There exists a transformer with $O(\log(1/\epsilon))$ layers and $O(K)$ heads such that the following holds: On any $(D; \mathbf{x}_{N+1})$ well-conditioned (cf. (5)) for all $f_k g_{k \in [K]}$, it outputs $\hat{\mathbf{w}}_{N+1} = \langle \hat{\mathbf{w}}; \mathbf{x}_{N+1} \rangle$, where*

$$\text{dist}(\hat{\mathbf{w}}; \text{conv} \{ \mathbf{w}_{\text{ridge}, \text{train}}^k : \mathbb{E}_{\text{val}}(\mathbf{w}_{\text{ridge}, \text{train}}^k) = \min_{k \in [K]} \mathbb{E}_{\text{val}}(\mathbf{w}_{\text{ridge}, \text{train}}^k) + g \}) \leq \epsilon$$

Above, $\mathbf{w}_{\text{ridge}, \text{train}}$ denotes the solution to (ICRidge) on the training split D_{train} .

4.1.1 Nearly Bayes-optimal ICL on noisy linear models with mixed noise levels

We build on Theorem 11 to show that transformers can perform nearly Bayes-optimal ICL when data come from noisy linear models with a mixture of K different noise levels $\sigma_1, \dots, \sigma_K > 0$.

Concretely, consider the following data generating model, where we first sample $\mathbf{P} = \mathbf{P}_{\mathbf{w}_{\mathcal{P}; k}}$ from $k \geq 2$ $([K])$, $\mathbf{w}_{\mathcal{P}} \sim \mathcal{N}(\mathbf{0}; \mathbf{I}_d)$, and then sample data $\{(\mathbf{x}_i; y_i)_{i \in [N+1]}\}_{\mathcal{P}} \stackrel{\text{iid}}{\sim} \mathbf{P}_{\mathbf{w}_{\mathcal{P}; k}}$ as

$$\mathbf{P}_{\mathbf{w}_{\mathcal{P}; k} : \mathbf{x}_i \sim \mathcal{N}(\mathbf{0}; \mathbf{I}_d); y_i = \langle \mathbf{x}_i; \mathbf{w}_{\mathcal{P}} \rangle + \mathcal{N}(0; \frac{\sigma_k^2}{k});$$

For any fixed $(N; d)$, consider the Bayes risk for predicting y_{N+1} under this model:

$$\text{BayesRisk} := \inf_A \mathbb{E} \frac{1}{2} (A(D)(\mathbf{x}_{N+1}) - y_{N+1})^2;$$

By standard Bayesian calculations, the above Bayes risk is attained when A is a certain mixture of K ridge regressions with regularization $\lambda_k = d \frac{\sigma_k^2}{N}$; however, the mixing weights depend on D in a highly non-trivial fashion (see Appendix J.2 for a derivation). By using the post-ICL validation mechanism in Theorem 11, we construct a transformer that achieves nearly the Bayes risk.

Theorem 12 (Nearly Bayes-optimal ICL; Informal version of Theorem J.1). *For sufficiently large $N; d$, there exists a transformer with $O(\log N)$ layers and $O(K)$ heads such that on the above model, it outputs a prediction \hat{y}_{N+1} that is nearly Bayes-optimal:*

$$\mathbb{E} \frac{1}{2} (y_{N+1} - \hat{y}_{N+1})^2 = \text{BayesRisk} + O((\log K=N)^{1=3}) \quad (8)$$

In particular, Theorem 12 applies in the *proportional setting* where $N; d$ are large and $N=d = (1) [22]$, in which case $\text{BayesRisk} = (1)$, and thus the transformer achieves vanishing excess risk relative to the Bayes risk at large N .

This substantially strengthens the results of Akyürek et al. [2], who empirically find that transformers can achieve nearly Bayes risk under any *fixed* noise level. By contrast, Theorem 12 shows that a *single* transformer can achieve nearly Bayes risk even under a mixture of K noise levels, with quantitative guarantees. Also, our proof in fact gives a stronger guarantee: The transformer approaches the *individual Bayes risks on all K noise levels simultaneously* (in addition to the overall Bayes risk for k as in Theorem 12). We demonstrate this empirically in Section 6 (cf. Figure 3b & 2).

Exact Bayes predictor vs. Post-ICL validation mechanism As BayesRisk is the theoretical lower bound for the risk of any possible ICL algorithm, Theorem 12 implies that our transformer performs similarly as the exact Bayes estimator⁵. Notice that our construction builds on the (generic) post-ICL validation mechanism, rather than a direct attempt of approximating the exact Bayes predictor, whose structure may vary significantly case-by-case. This highlights post-ICL validation as a promising mechanism for approximating the Bayes predictor on broader classes of problems beyond noisy linear models, which we leave as future work.

Generalized linear models with adaptive link function selection As another example of the post-ICL validation mechanism, we construct a transformer that can learn a generalized linear model with adaptively chosen link function for the particular ICL instance; see Theorem J.2.

4.2 Pre-ICL testing mechanism

In our second mechanism, pre-ICL testing, the transformer runs a *distribution testing* procedure on the input sequence to determine the right ICL algorithm to use. While the test (and thus the mechanism itself) could in principle be general, we focus on cases where the test amounts to computing some simple summary statistics of the input sequence.

To showcase pre-ICL testing, we consider the toy problem of selecting between in-context regression and in-context classification, by running the following *binary type check* on the input labels $\{y_i\}_{i \in [N]}$.

$$\text{binary}(D) = \frac{1}{N} \sum_{i=1}^N (y_i); \quad (y) := \begin{cases} \frac{y}{\sigma} < 1; & y \geq \tau_0; 1g; \\ 0; & y \notin ["; "] [[1"; 1 + "; \\ \text{linear interpolation}; & \text{otherwise}; \end{cases}$$

⁵By the Bayes risk decomposition for square loss, (8) implies that $\mathbb{E}[(\hat{y}_{N+1} - y_{N+1})^2] = O((\log K=N)^{1=3})$.

Lemma 13. *There exists a single attention layer with 6 heads that implements $\text{argmax}_{y \in \mathcal{Y}} \sum_{i=1}^n \mathbb{1}\{y = y_i\}$ exactly.*

Using this test, we construct a transformer that performs logistic regression when labels are binary, and linear regression with high probability if the label admits a continuous distribution.

Proposition 14 (Adaptive regression or classification; Informal version of Proposition I.4). *There exists a transformer with $O(\log(1/\epsilon))$ layers such that the following holds: On any D such that $y_i \geq f_0 + 1/g$, it outputs \hat{y}_{N+1} that ϵ -approximates the prediction of in-context logistic regression.*

By contrast, for any distribution P whose marginal distribution of y is not concentrated around $f_0 + 1/g$, with high probability (over D), \hat{y}_{N+1} ϵ -approximates the prediction of in-context least squares.

The proofs can be found in Appendix I.3. We additionally show that transformers can implement more complex tests such as a *linear correlation test*, which can be useful in certain scenarios such as “confident linear regression” (predict only when the signal-to-noise ratio is high); see Appendix I.4.

5 Analysis of pretraining

Building on the expressivity results in Section 3 & 4, we provide the first line of polynomial sample complexity results for *pretraining* transformers to perform ICL (including with in-context algorithm selection). We begin by providing a generic generalization guarantee for pretraining transformers.

Consider the pretraining ERM problem (TF-ERM), which minimizes the pretraining risk $\hat{L}_{\text{icl}}(\cdot)$ over n pretraining sequences. Let $L_{\text{icl}}(\cdot)$ denote the corresponding population risk.

Theorem 15 (Generalization of transformers; Informal version of Theorem K.1). *The solution \hat{b} to (TF-ERM) over transformers with L layers, M heads per layer, and hidden dimension D' satisfies*

$$L_{\text{icl}}(\hat{b}) \leq \inf L_{\text{icl}}(\cdot) + \Theta\left(\frac{\sqrt{L^2(MD^2 + DD')}}{n}\right);$$

Theorem 15 builds on standard uniform concentration analysis via chaining (Proposition B.4). Combining Theorem 15 with the in-context linear regression construction in Theorem 4 gives the following end-to-end result on the excess in-context prediction risk of trained transformers.

Theorem 16 (Pretraining transformers for in-context linear regression; Informal version of Theorem K.2). *Under Assumption A and $N = \Theta(d)$, the solution \hat{b} to (TF-ERM) with $L = O(\log(N/\epsilon))$ layers, $M = 3$ heads, $D' = 0$ (attention-only as in Theorem 4) achieves small excess ICL risk over the best linear predictor $\mathbf{w}_P^* := \mathbb{E}_P[\mathbf{x}\mathbf{x}^\top]^{-1}\mathbb{E}_P[\mathbf{x}y]$ for each P :*

$$L_{\text{icl}}(\hat{b}) - \mathbb{E}_{P \sim \mathcal{P}} \mathbb{E}_{(\mathbf{x}; y) \sim P} \frac{1}{2} (y - \langle \mathbf{w}_P^*, \mathbf{x} \rangle)^2 \leq \Theta\left(\frac{d^2}{n} + \frac{d^2}{N}\right);$$

See Appendix K.2 for similar results in several additional settings.

6 Experiments

We test our theory by studying the ICL and in-context algorithm selection capabilities of transformers, using the encoder-based architecture in our theoretical constructions (Definition 3). Due to limited space, additional experimental details can be found in Appendix M.1. Results with a decoder architecture as in [31, 47] (including the setup of Figure 2) can be found in Appendix M.2.

Training data distributions and evaluation We train a 12-layer transformer, with two modes for the training sequence (instance) distribution \mathcal{P} . In the “base” mode, similar to [31, 2, 86, 47], we sample the training instances from *one* of the following base distributions (tasks), where we first sample $P = P_{\mathbf{w}_?}$ by sampling $\mathbf{w}_? \sim \mathcal{N}(\mathbf{0}; I_{d=d})$, and then sample $f(\mathbf{x}_i; y_i) g_{i \in [N+1]} \stackrel{\text{iid}}{\sim} P_{\mathbf{w}_?}$ as $\mathbf{x}_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}; I_d)$, and y_i from one of the following models studied in Section 3:

1. Linear model: $y_i = \langle \mathbf{w}_?, \mathbf{x}_i \rangle$;
2. Noisy linear model: $y_i = \langle \mathbf{w}_?, \mathbf{x}_i \rangle + z_i$, where $\sigma > 0$ is a fixed noise level, and $z_i \sim \mathcal{N}(0; 1)$.

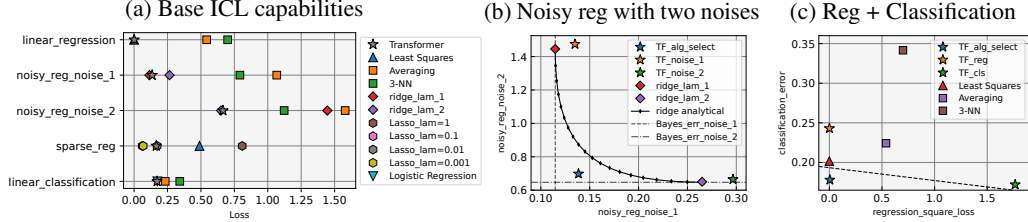


Figure 3: ICL capabilities of the transformer architecture used in our theoretical constructions. (a) On five representative base tasks, transformers approximately match the best baseline algorithm for each task, when pretrained on the corresponding task. (b,c) A **single transformer** `TF_alg_select` **simultaneously approaches the performance of the strongest baseline algorithm** on two separate tasks: (b) noisy linear regression with two different noise levels $\sim \mathcal{N}(0,1;0.5g)$, and (c) adaptively selecting between regression and classification.

3. Sparse linear model: $y_i = \langle \mathbf{w}_\gamma, \mathbf{x}_i \rangle$ with $\|\mathbf{w}_\gamma\|_0 = s$, where $s < d$ is a fixed sparsity level, and in this case we sample \mathbf{w}_γ from a special prior supported on s -sparse vectors;
4. Linear classification model: $y_i = \text{sign}(\langle \mathbf{w}_\gamma, \mathbf{x}_i \rangle)$.

These base tasks have been empirically investigated by Garg et al. [31], though we remark that our architecture (used in our theory) differs from theirs in several aspects, such as encoder-based architecture instead of decoder-based, and ReLU activation instead of softmax. All experiments use $d = 20$. We choose $\mathcal{N}(0,1;0.5g)$ and $N = 20$ for noisy linear regression, $s = 3$ and $N = 10$ for sparse linear regression, and $N = 40$ for linear regression and linear classification.

In the “mixture” mode, \mathcal{D} is the uniform *mixture of two or more base distributions*. We consider two representative mixture modes studied in Section 4:

- Linear model + linear classification model;
- Noisy linear model with four noise levels $\sim \mathcal{N}(0,1;0.25;0.5;1g)$.

Transformers trained with the mixture mode will be evaluated on *multiple* base distributions simultaneously. When the base distributions are sufficiently diverse, a transformer performing well on all of them will *likely* be performing some level of in-context algorithm selection. We evaluate transformers against standard machine learning algorithms in context (for each task respectively) as baselines.

Results Figure 3a shows the ICL performance of transformers on five base tasks, within each the transformer is trained on the same task. Transformers match the best baseline algorithm in four out of the five cases, except for the sparse regression task where the Transformer still outperforms least squares and matches Lasso with some choices of λ (thus utilizing sparsity to some extent). This demonstrates the strong ICL capability of the transformer architecture considered in our theory.

Figure 3b & 3c examine the in-context algorithm selection capability of transformers, on noisy linear regression with two different noise levels (Figure 3b), and regression + classification (Figure 3c). In both figures, the transformer trained in the mixture mode (`TF_alg_select`) approaches the best baseline algorithm on both tasks simultaneously. By contrast, transformers trained in the base mode for one of the tasks perform well on that task but behave suboptimally on the other task as expected. The existence of `TF_alg_select` showcases a single transformer that performs well on multiple tasks simultaneously (and thus has to perform in-context algorithm selection to some extent), supporting our theoretical results in Section 4.

7 Conclusion

This work shows that transformers can perform complex in-context learning procedures with strong in-context algorithm selection capabilities, by both explicit theoretical constructions and experiments. We believe our work opens up many exciting directions, such as (1) more mechanisms for in-context algorithm selection; (2) Bayes-optimal ICL on other problems by either the post-ICL validation mechanism or new approaches; (3) understanding the internal workings of transformers performing in-context algorithm selection; (4) other mechanisms for implementing complex ICL procedures beyond in-context algorithm selection; (5) further statistical analyses, e.g. of pretraining. Besides, this work focuses on the transformer architecture; alternative sequence-to-sequence architectures (such as RNNs) are beyond our scope but would be interesting directions for future work.

Acknowledgment

The authors would like to thank Tengyu Ma and Jason D. Lee for the many insightful discussions. S. Mei is supported in part by NSF DMS-2210827 and NSF CCF-2315725.

References

- [1] A. Agarwal, S. Negahban, and M. J. Wainwright. Fast global convergence rates of gradient methods for high-dimensional statistical recovery. *Advances in Neural Information Processing Systems*, 23, 2010.
- [2] E. Akyürek, D. Schuurmans, J. Andreas, T. Ma, and D. Zhou. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- [3] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [4] F. Bach. Breaking the curse of dimensionality with convex neural networks. *The Journal of Machine Learning Research*, 18(1):629–681, 2017.
- [5] Y. Bai, M. Chen, P. Zhou, T. Zhao, J. Lee, S. Kakade, H. Wang, and C. Xiong. How important is the train-validation split in meta-learning? In *International Conference on Machine Learning*, pages 543–553. PMLR, 2021.
- [6] Y. Bai, S. Mei, H. Wang, and C. Xiong. Don’t just blame over-parametrization for over-confidence: Theoretical analysis of calibration in binary classification. In *International Conference on Machine Learning*, pages 566–576. PMLR, 2021.
- [7] J. Baxter. A model of inductive bias learning. *Journal of artificial intelligence research*, 12: 149–198, 2000.
- [8] A. Beck and M. Teboulle. Gradient-based algorithms with applications to signal recovery. *Convex optimization in signal processing and communications*, pages 42–88, 2009.
- [9] S. Bengio, Y. Bengio, J. Cloutier, and J. Gescei. On the optimization of a synaptic learning rule. In *Optimality in Biological and Artificial Networks?*, pages 281–303. Routledge, 2013.
- [10] S. Bhattamishra, K. Ahuja, and N. Goyal. On the ability and limitations of transformers to recognize formal languages. *arXiv preprint arXiv:2009.11264*, 2020.
- [11] S. Bhattamishra, A. Patel, and N. Goyal. On the computational power of transformers and its implications in sequence modeling. *arXiv preprint arXiv:2006.09286*, 2020.
- [12] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [13] S. Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- [14] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [15] S. Chan, A. Santoro, A. Lampinen, J. Wang, A. Singh, P. Richemond, J. McClelland, and F. Hill. Data distributional properties drive emergent in-context learning in transformers. *Advances in Neural Information Processing Systems*, 35:18878–18891, 2022.
- [16] L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.

- [17] K. Chua, Q. Lei, and J. D. Lee. How fine-tuning allows for effective meta-learning. *Advances in Neural Information Processing Systems*, 34:8871–8884, 2021.
- [18] D. Dai, Y. Sun, L. Dong, Y. Hao, Z. Sui, and F. Wei. Why can gpt learn in-context? language models secretly perform gradient descent as meta optimizers. *arXiv preprint arXiv:2212.10559*, 2022.
- [19] G. Denevi, C. Ciliberto, D. Stamos, and M. Pontil. Incremental learning-to-learn with statistical guarantees. *arXiv preprint arXiv:1803.08089*, 2018.
- [20] G. Denevi, C. Ciliberto, D. Stamos, and M. Pontil. Learning to learn around a common mean. *Advances in Neural Information Processing Systems*, 31, 2018.
- [21] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [22] E. Dobriban and S. Wager. High-dimensional asymptotics of prediction: Ridge regression and classification. *The Annals of Statistics*, 46(1):247–279, 2018.
- [23] L. Dong, S. Xu, and B. Xu. Speech-transformer: a no-recurrence sequence-to-sequence model for speech recognition. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5884–5888. IEEE, 2018.
- [24] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, and Z. Sui. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- [25] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [26] S. S. Du, W. Hu, S. M. Kakade, J. D. Lee, and Q. Lei. Few-shot learning via learning the representation, provably. *arXiv preprint arXiv:2002.09434*, 2020.
- [27] B. L. Edelman, S. Goel, S. Kakade, and C. Zhang. Inductive biases and variable creation in self-attention mechanisms. In *International Conference on Machine Learning*, pages 5793–5831. PMLR, 2022.
- [28] N. Elhage, N. Nanda, C. Olsson, T. Henighan, N. Joseph, B. Mann, A. Askell, Y. Bai, A. Chen, T. Conerly, et al. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021.
- [29] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017.
- [30] C. Finn, A. Rajeswaran, S. Kakade, and S. Levine. Online meta-learning. In *International Conference on Machine Learning*, pages 1920–1930. PMLR, 2019.
- [31] S. Garg, D. Tsipras, P. S. Liang, and G. Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- [32] A. Giannou, S. Rajput, J.-y. Sohn, K. Lee, J. D. Lee, and D. Papailiopoulos. Looped transformers as programmable computers. *arXiv preprint arXiv:2301.13196*, 2023.
- [33] M. Hahn. Theoretical limitations of self-attention in neural sequence models. *Transactions of the Association for Computational Linguistics*, 8:156–171, 2020.
- [34] S. Hochreiter, A. S. Younger, and P. R. Conwell. Learning to learn using gradient descent. In *Artificial Neural Networks—ICANN 2001: International Conference Vienna, Austria, August 21–25, 2001 Proceedings 11*, pages 87–94. Springer, 2001.
- [35] N. Hollmann, S. Müller, K. Eggenberger, and F. Hutter. TabPFN: A transformer that solves small tabular classification problems in a second. *arXiv preprint arXiv:2207.01848*, 2022.

- [36] T. Hospedales, A. Antoniou, P. Micaelli, and A. Storkey. Meta-learning in neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(9):5149–5169, 2021.
- [37] J. Hron, Y. Bahri, J. Sohl-Dickstein, and R. Novak. Infinite attention: Nngp and ntk for deep attention networks. In *International Conference on Machine Learning*, pages 4376–4386. PMLR, 2020.
- [38] D. Hsu, S. M. Kakade, and T. Zhang. Random design analysis of ridge regression. In *Conference on learning theory*, pages 9–1. JMLR Workshop and Conference Proceedings, 2012.
- [39] S. Jelassi, M. E. Sander, and Y. Li. Vision transformers provably learn spatial structure. *arXiv preprint arXiv:2210.09221*, 2022.
- [40] K. Ji, J. D. Lee, Y. Liang, and H. V. Poor. Convergence of meta-learning with task-specific adaptation over partial parameters. *Advances in Neural Information Processing Systems*, 33: 11490–11500, 2020.
- [41] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. vZidek, A. Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [42] S. M. Kakade, V. Kanade, O. Shamir, and A. Kalai. Efficient learning of generalized linear and single index models with isotonic regression. *Advances in Neural Information Processing Systems*, 24, 2011.
- [43] M. Khodak, M.-F. F. Balcan, and A. S. Talwalkar. Adaptive gradient-based meta-learning methods. *Advances in Neural Information Processing Systems*, 32, 2019.
- [44] L. Kirsch and J. Schmidhuber. Meta learning backpropagation and improving it. *Advances in Neural Information Processing Systems*, 34:14122–14134, 2021.
- [45] L. Kirsch, J. Harrison, J. Sohl-Dickstein, and L. Metz. General-purpose in-context learning by meta-learning transformers. *arXiv preprint arXiv:2212.04458*, 2022.
- [46] K. Li and J. Malik. Learning to optimize. *arXiv preprint arXiv:1606.01885*, 2016.
- [47] Y. Li, M. E. Ildiz, D. Papailiopoulos, and S. Oymak. Transformers as algorithms: Generalization and implicit model selection in in-context learning. *arXiv preprint arXiv:2301.07067*, 2023.
- [48] B. Liu, J. T. Ash, S. Goel, A. Krishnamurthy, and C. Zhang. Transformers learn shortcuts to automata. *arXiv preprint arXiv:2210.10749*, 2022.
- [49] J. Liu, D. Shen, Y. Zhang, B. Dolan, L. Carin, and W. Chen. What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*, 2021.
- [50] Y. Lu, M. Bartolo, A. Moore, S. Riedel, and P. Stenetorp. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv preprint arXiv:2104.08786*, 2021.
- [51] A. Madani, B. McCann, N. Naik, N. S. Keskar, N. Anand, R. R. Eguchi, P.-S. Huang, and R. Socher. Progen: Language modeling for protein generation. *arXiv preprint arXiv:2004.03497*, 2020.
- [52] A. Maurer, M. Pontil, and B. Romera-Paredes. The benefit of multitask representation learning. *Journal of Machine Learning Research*, 17(81):1–32, 2016.
- [53] P. McCullagh. *Generalized linear models*. Routledge, 2019.
- [54] S. Mei, Y. Bai, and A. Montanari. The landscape of empirical risk for nonconvex losses. *The Annals of Statistics*, 46(6A):2747–2774, 2018.

- [55] S. Min, M. Lewis, H. Hajishirzi, and L. Zettlemoyer. Noisy channel language model prompting for few-shot text classification. *arXiv preprint arXiv:2108.04106*, 2021.
- [56] S. Min, M. Lewis, L. Zettlemoyer, and H. Hajishirzi. Metaicl: Learning to learn in context. *arXiv preprint arXiv:2110.15943*, 2021.
- [57] S. Min, X. Lyu, A. Holtzman, M. Artetxe, M. Lewis, H. Hajishirzi, and L. Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*, 2022.
- [58] N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel. A simple neural attentive meta-learner. *arXiv preprint arXiv:1707.03141*, 2017.
- [59] S. Müller, N. Hollmann, S. P. Arango, J. Grabocka, and F. Hutter. Transformers can do bayesian inference. *arXiv preprint arXiv:2112.10510*, 2021.
- [60] T. Nagler. Statistical foundations of prior-data fitted networks. *arXiv preprint arXiv:2305.11097*, 2023.
- [61] D. K. Naik and R. J. Mammone. Meta-neural networks that learn by learning. In *[Proceedings 1992] IJCNN International Joint Conference on Neural Networks*, volume 1, pages 437–442. IEEE, 1992.
- [62] S. N. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu. A unified framework for high-dimensional analysis of m-estimators with decomposable regularizers. 2012.
- [63] Y. Nesterov. *Lectures on convex optimization*, volume 137. Springer, 2018.
- [64] C. Olsson, N. Elhage, N. Nanda, N. Joseph, N. DasSarma, T. Henighan, B. Mann, A. Askell, Y. Bai, A. Chen, et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- [65] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [66] N. Parikh, S. Boyd, et al. Proximal algorithms. *Foundations and trends® in Optimization*, 1(3):127–239, 2014.
- [67] J. Pérez, J. Marinković, and P. Barceló. On the turing completeness of modern neural network architectures. *arXiv preprint arXiv:1901.03429*, 2019.
- [68] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [69] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [70] A. Raventos, M. Paul, F. Chen, and S. Ganguli. The effects of pretraining task diversity on in-context learning of ridge regression. In *ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2023.
- [71] S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. In *International conference on learning representations*, 2017.
- [72] Y. Razeghi, R. L. Logan IV, M. Gardner, and S. Singh. Impact of pretraining term frequencies on few-shot reasoning. *arXiv preprint arXiv:2202.07206*, 2022.
- [73] S. Reed, K. Zolna, E. Parisotto, S. G. Colmenarejo, A. Novikov, G. Barth-Maron, M. Gimenez, Y. Sulsky, J. Kay, J. T. Springenberg, et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022.
- [74] O. Rubin, J. Herzig, and J. Berant. Learning to retrieve prompts for in-context learning. *arXiv preprint arXiv:2112.08633*, 2021.

- [75] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap. Meta-learning with memory-augmented neural networks. In *International conference on machine learning*, pages 1842–1850. PMLR, 2016.
- [76] N. Saunshi, A. Gupta, and W. Hu. A representation learning perspective on the importance of train-validation splitting in meta-learning. In *International Conference on Machine Learning*, pages 9333–9343. PMLR, 2021.
- [77] J. Schmidhuber. *Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook*. PhD thesis, Technische Universität München, 1987.
- [78] K. Shen, J. Guo, X. Tan, S. Tang, R. Wang, and J. Bian. A study on relu and softmax in transformer. *arXiv preprint arXiv:2302.06461*, 2023.
- [79] C. Snell, R. Zhong, D. Klein, and J. Steinhardt. Approximating how single head attention learns. *arXiv preprint arXiv:2103.07601*, 2021.
- [80] J. Snell, K. Swersky, and R. Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017.
- [81] S. Thrun and L. Pratt. *Learning to learn*. Springer Science & Business Media, 2012.
- [82] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- [83] N. Tripuraneni, M. Jordan, and C. Jin. On the theory of transfer learning: The importance of task diversity. *Advances in neural information processing systems*, 33:7852–7862, 2020.
- [84] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [85] R. Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [86] J. von Oswald, E. Niklasson, E. Randazzo, J. Sacramento, A. Mordvintsev, A. Zhmoginov, and M. Vladymyrov. Transformers learn in-context by gradient descent. *arXiv preprint arXiv:2212.07677*, 2022.
- [87] M. J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- [88] X. Wang, S. Yuan, C. Wu, and R. Ge. Guarantees for tuning the step size using a learning-to-learn approach. In *International Conference on Machine Learning*, pages 10981–10990. PMLR, 2021.
- [89] C. Wei, Y. Chen, and T. Ma. Statistically meaningful approximation: a case study on approximating turing machines with transformers. *arXiv preprint arXiv:2107.13163*, 2021.
- [90] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- [91] J. Wei, J. Wei, Y. Tay, D. Tran, A. Webson, Y. Lu, X. Chen, H. Liu, D. Huang, D. Zhou, et al. Larger language models do in-context learning differently. *arXiv preprint arXiv:2303.03846*, 2023.
- [92] G. Weiss, Y. Goldberg, and E. Yahav. Thinking like transformers. In *International Conference on Machine Learning*, pages 11080–11090. PMLR, 2021.
- [93] M. Wortsman, J. Lee, J. Gilmer, and S. Kornblith. Replacing softmax with relu in vision transformers. *arXiv preprint arXiv:2309.08586*, 2023.
- [94] S. M. Xie, A. Raghunathan, P. Liang, and T. Ma. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.

- [95] S. Yao, B. Peng, C. Papadimitriou, and K. Narasimhan. Self-attention networks can process bounded hierarchical languages. *arXiv preprint arXiv:2105.11115*, 2021.
- [96] C. Ying, T. Cai, S. Luo, S. Zheng, G. Ke, D. He, Y. Shen, and T.-Y. Liu. Do transformers really perform badly for graph representation? *Advances in Neural Information Processing Systems*, 34:28877–28888, 2021.
- [97] C. Yun, S. Bhojanapalli, A. S. Rawat, S. J. Reddi, and S. Kumar. Are transformers universal approximators of sequence-to-sequence functions? *arXiv preprint arXiv:1912.10077*, 2019.
- [98] Y. Zhang, A. Backurs, S. Bubeck, R. Eldan, S. Gunasekar, and T. Wagner. Unveiling transformers with lego: a synthetic reasoning task. *arXiv preprint arXiv:2206.04301*, 2022.
- [99] Y. Zhang, B. Liu, Q. Cai, L. Wang, and Z. Wang. An analysis of attention via the lens of exchangeability and latent variable models. *arXiv preprint arXiv:2212.14852*, 2022.
- [100] Z. Zhao, E. Wallace, S. Feng, D. Klein, and S. Singh. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*, pages 12697–12706. PMLR, 2021.
- [101] X. Zuo, Z. Chen, H. Yao, Y. Cao, and Q. Gu. Understanding train-validation split in meta-learning with neural networks. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=JVlyfHEEmOk>.

A Related work

In-context learning The in-context learning (ICL) capability of large language models (LLMs) has gained significant attention since demonstrated on GPT-3 Brown et al. [12]. A number of subsequent empirical studies have contributed to a better understanding of the capabilities and limitations of ICL in LLM systems, which include but are not limited to [49, 55, 56, 50, 100, 74, 72, 28, 45, 91]. For an overview of ICL, see the survey by Dong et al. [24] which highlights some key findings and advancements in this direction.

A line of recent work investigates why and how LLMs perform ICL [94, 31, 86, 2, 18, 32, 47, 70]. In particular, Xie et al. [94] propose a Bayesian inference framework explaining how ICL works despite formatting differences between training and inference distributions. Garg et al. [31] show empirically that transformers could be trained from scratch to perform ICL of linear models, sparse linear models, two-layer neural networks, and decision trees. Li et al. [47] analyze the generalization error of trained ICL transformers from a stability viewpoint. They also experimentally show that transformers could perform “in-context model selection” (conceptually similar to in-context algorithm selection considered in this work) in specific tasks and presented related theoretical hypotheses. However, they do not provide concrete mechanisms or constructions for in-context model selection. A recent work [99] shows that pretrained transformers can perform Bayesian inference in latent variable models, which may also be interpreted as a mechanism for ICL. Our experimental findings extend these results by unveiling and demonstrating the in-context algorithm selection capabilities of transformers.

Closely related to our theoretical results are [86, 2, 18, 32], which show (among many things) that transformers can perform ICL by simulating gradient descent. However, these results do not provide quantitative error bounds for simulating multi-step gradient descent, and only handle linear regression models or their simple variants. Among these works, Akyürek et al. [2] showed that transformers can implement learning algorithms for linear models based on gradient descent and closed-form ridge regression; it also presented preliminary evidence that learned transformers perform ICL similar to Bayes-optimal ridge regression. Our work builds upon and substantially extends this line of work by (1) providing a more efficient construction for in-context gradient descent; (2) providing an end-to-end theory with additional results for pretraining and statistical power; (3) analyzing a broader spectrum of ICL algorithms, including least squares, ridge regression, Lasso, convex risk minimization for generalized linear models, and gradient descent on two-layer neural networks; and (4) constructing more complex ICL procedures using in-context algorithm selection.

When in-context data are generated from a prior, the Bayes risk is a theoretical lower bound for the risk of any possible ICL algorithm, including transformers. Xie et al. [94], Akyürek et al. [2] observe

that learned transformers behave closely to the Bayes predictor on a variety of tasks such as hidden Markov models [94] and noisy linear regression with a fixed noise level [2, 47]. Using the in-context algorithm selection mechanism (more precisely the post-ICL validation mechanism), we show that transformers can perform nearly-Bayes optimal ICL in noisy linear models with mixed noise levels (a strictly more challenging task than considered in [2, 47]), with both concrete theoretical guarantees (Section 4.1.1) and empirical evidence (Figure 2 & 3b). Complementary to these works, a line of work on “prior-data fitted networks” [59, 60, 35] also empirically demonstrates the Bayesian optimality of transformers in various settings. Our expressivity results support these empirical findings and are applicable beyond the Bayesian setting, e.g. for providing frequentist in-context prediction guarantees for transformers.

Transformers and its theory The transformer architecture, introduced by [84], has revolutionized natural language processing and been adopted in most of the recently developed large language models such as BERT and GPT [68, 21, 12]. Broader, transformers have demonstrated remarkable performance in many other fields of artificial intelligence such as computer vision, speech, graph processing, reinforcement learning, and biological applications [23, 25, 51, 69, 96, 16, 41, 73, 65, 14]. Towards a better theoretical understanding, recent work has studied the capabilities [97, 67, 37, 95, 11, 98, 48], limitations [33, 10], and internal workings [28, 79, 92, 27, 64] of transformers.

We remark that the transformer architecture used in our theoretical constructions differs from the standard one by replacing the softmax activation (in the attention layers) with a (normalized) ReLU function. Transformers with ReLU activations is experimentally studied in the recent work of Shen et al. [78], who find that they perform as well as the standard softmax activation in many NLP tasks.

Meta-learning Training models (such as transformers) to perform ICL can be viewed as an approach for the broader problem of learning-to-learn or meta-learning [77, 61, 81]. A number of other approaches has been studied extensively for this problem, including (and not limited to) training a meta-learner on how to update the parameters of a downstream learner [9, 46], learning parameter initializations that quickly adapt to downstream tasks [29, 71], learning latent embeddings that allow for effective similarity search [80]. Most relevant to the ICL setting are approaches that directly take as input examples from a downstream task and a query input and produce the corresponding output [34, 58, 75, 44]. For a comprehensive overview, see the survey [36].

Theoretical aspects of meta-learning have received significant recent interest [7, 52, 26, 83, 19, 30, 43, 40, 88, 20, 5, 76, 17, 101]. In particular, [52, 26, 83] analyzed the benefit of multi-task learning through a representation learning perspective, and [88, 20, 5, 76, 101] studied the statistical properties of learning the parameter initialization for downstream tasks.

Techniques We build on various existing techniques from the statistics and learning theory literature to establish our approximation and generalization guarantees for transformers. For the approximation component, we rely on a technical result of Bach [4] on the approximation power of ReLU networks. We use this result to show that transformers can approximate gradient descent (GD) on a broad range of loss functions, substantially extending the results of [86, 2, 18] who primarily consider the square loss. The recent work of Giannou et al. [32] also approximates GD with general loss functions by transformers, though using a different technique of forcing the softmax activations to act as sigmoids. Our analyses of Lasso and generalized linear models build on [87, 62, 1, 54]. Our generalization bound for transformers (used in our pretraining results) build on a chaining argument [87].

B Technical tools

Additional notation for proofs We say a random variable X is σ^2 -sub-Gaussian (or $SG(\sigma^2)$ interchangeably) if $E[\exp(X^2/\sigma^2)] \leq 2$. A random vector $\mathbf{x} \in \mathbb{R}^d$ is σ^2 -sub-Gaussian if $\langle \mathbf{v}, \mathbf{x} \rangle$ is σ^2 -sub-Gaussian for all $\|\mathbf{v}\|_2 = 1$. A random variable X is K -sub-Exponential (or $SE(K)$ interchangeably) if $E[\exp(X/K)] \leq 2$.

B.1 Concentration inequalities

Lemma B.1. Let $\mathbf{x} \sim N(\mathbf{0}; \mathbf{I}_{d=d})$. Then we have

$$P\left(\|\mathbf{x}\|_2 \geq \sqrt{d} + t\right) \leq e^{-d^{-1}t^2}.$$

Lemma B.2 (Theorem 6.1 of [87]). Let $X = [X_{ij}] \in \mathbb{R}^{n \times d}$ be a Gaussian random matrix with $X_{ij} \sim \mathcal{N}(0; 1)$. Let $\lambda_{\min}(X)$ and $\lambda_{\max}(X)$ be the minimum and maximum singular value of X , respectively. Then we have

$$\begin{aligned} \mathbb{P} \left[\lambda_{\max}(X) \leq \sqrt{\frac{d}{n}} + 1 + \sqrt{\frac{d}{n}} e^{-n^{-2/2}} \right]; \\ \mathbb{P} \left[\lambda_{\min}(X) \geq \sqrt{\frac{d}{n}} - 1 - \sqrt{\frac{d}{n}} e^{-n^{-2/2}} \right]; \end{aligned}$$

The following lemma is a standard result of covariance concentration, see e.g. [85, Theorem 4.6.1].

Lemma B.3. Suppose that $\mathbf{x}_1, \dots, \mathbf{x}_N$ are independent d -dimensional K -sub-Gaussian random vectors. Then as long as $N \geq C_0 d$, with probability at least $1 - \exp(-N/C_0)$ we have

$$\left| \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top - \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top] \right|_{\text{op}} \leq 8K^2;$$

where C_0 is a universal constant.

Lemma B.4. For random matrix $\mathbf{X} = [X_{ij}] \in \mathbb{R}^{N \times d}$ with $X_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0; 1)$ and $\mathbf{u} = [u_j] \in \mathbb{R}^N$ with $u_j \stackrel{\text{iid}}{\sim} \mathcal{N}(0; \sigma^2)$, it holds that

$$\mathbb{P} \left[\|\mathbf{X}^\top \mathbf{u}\|_\infty \geq \sqrt{\frac{d}{N} \log(2d)} + \exp(-N/2) \right];$$

Proof. We consider $\mathbf{u}_j := [X_{ij}]_i \in \mathbb{R}^N$, then $\|\mathbf{X}^\top \mathbf{u}\|_\infty = \max_{j \in [d]} \|\mathbf{u}_j\|_2$. Notice that the random variables $\|\mathbf{u}_1\|_2, \dots, \|\mathbf{u}_d\|_2$ are independent $\mathcal{N}(0; k^2 \sigma^2)$, and hence

$$\mathbb{P} \left[\max_{j \in [d]} \|\mathbf{u}_j\|_2 \geq t \right] \leq 2d \exp\left(-\frac{t^2}{2k^2 \sigma^2}\right);$$

Further, by Lemma B.1, $\mathbb{P}(k^2 \sigma^2 \geq 2\sqrt{N}) \leq \exp(-N/2)$. Taking $t = \sqrt{\frac{d}{8N} \log(2d)}$ completes the proof. \square

B.2 Approximation theory

For any signed measure ν over a space W , let $\text{TV}(\nu) := \int_W |\nu| d\mathbf{w}$ denote its total measure. Recall $\sigma(\cdot) = \text{ReLU}(\cdot)$ is the standard relu activation, and $B_\infty^k(R) = [-R; R]^k$ denotes the standard ℓ_∞ ball in \mathbb{R}^k with radius $R > 0$.

Definition B.1 (Sufficiently smooth k -variable function). We say a function $g: \mathbb{R}^k \rightarrow \mathbb{R}$ is $(R; C)$ -smooth, if for $s = d(k-1) + 2e + 2$, g is a C^s function on $B_\infty^k(R)$, and

$$\sup_{\mathbf{z} \in B_\infty^k(R)} \|\mathbf{r}^i g(\mathbf{z})\|_\infty = \sup_{\mathbf{z} \in B_\infty^k(R)} \max_{j_1, \dots, j_i \in [k]} |g(\mathbf{x})|_{j_1, \dots, j_i} \leq L_i$$

for all $i \in \{0, 1, \dots, s\}$, with $\max_{0 \leq i \leq s} L_i R^i \leq C$.

The following result for expressing smooth functions as a random feature model with relu activation is adapted from Bach [4, Proposition 5].

Lemma B.5 (Expressing sufficiently smooth functions by relu random features). Suppose function $g: \mathbb{R}^k \rightarrow \mathbb{R}$ is $(R; C)$ -smooth. Then there exists a signed measure ν over $W = \mathbb{R}^{k+1}$ with $\|\nu\|_1 = 1$ such that

$$g(\mathbf{x}) = \int_W \frac{1}{R} (\mathbf{w}^\top [\mathbf{x}; R]) d\nu(\mathbf{w}); \quad \forall \mathbf{x} \in \mathbb{R}^k$$

and $\text{TV}(\nu) \leq C(k)C$, where $C(k) < 1$ is a constant that only depends on k .

Lemma B.6 (Uniform finite-neuron approximation). Let X be a space equipped with a distance function $d_X(\cdot, \cdot): X \times X \rightarrow \mathbb{R}_{\geq 0}$. Suppose function $g: X \rightarrow \mathbb{R}$ is given by

$$g(\mathbf{x}) = \int_W (\mathbf{x}; \mathbf{w}) d\nu(\mathbf{w});$$

where $(\cdot; \cdot) : X \rightarrow \mathbb{R}$ is L -Lipschitz in d_X in the first argument, and $\int_W (\cdot; \cdot) d\mathbf{w}$ is a signed measure over W with finite total measure $A = \text{TV}(\cdot; \cdot) < 1$. Then for any $\epsilon > 0$, there exists $K = O(A^2 B^2 \log N(X; d_X; \frac{\epsilon}{3AL}))$, such that

$$\sup_{\mathbf{x} \in X} g(\mathbf{x}) = \frac{A}{K} \sum_{i=1}^K (\mathbf{x}; \mathbf{w}_i) \quad \epsilon;$$

where $N(X; d_X; \frac{\epsilon}{3AL})$ denotes the $(\frac{\epsilon}{3AL})$ -covering number of X in d_X .

Proof. Let $(\mathbf{w}) := \text{sign}(\int_W (\cdot; \mathbf{w}) d\mathbf{w})$ denote the sign of the density $d(\mathbf{w})$. We have

$$g(\mathbf{x}) = A \int_W (\mathbf{w}) (\mathbf{x}; \mathbf{w}) \frac{jd(\mathbf{w})}{A} d\mathbf{w}. \quad (9)$$

Note that $jd(\mathbf{w})/A$ is the density of a probability distribution over W . Thus for any $\mathbf{x} \in X$, as long as $K = O(A^2 B^2 \log(1/\epsilon))$, we can sample $\mathbf{w}_1, \dots, \mathbf{w}_K \stackrel{\text{iid}}{\sim} jd(\cdot)/A$, and obtain by Hoeffding's inequality that with probability at least $1 - \epsilon$,

$$g(\mathbf{x}) = \frac{A}{K} \sum_{i=1}^K (\mathbf{w}_i) (\mathbf{x}; \mathbf{w}_i) \quad \epsilon;$$

Let $N(\frac{\epsilon}{3AL}) := N(X; d_X; \frac{\epsilon}{3AL})$ for shorthand. By union bound, as long as $K = O(A^2 B^2 \log(N(\frac{\epsilon}{3AL})))$, we have with probability at least $1 - \epsilon$ that for every \mathbf{x} in the covering set corresponding to $N(\frac{\epsilon}{3AL})$,

$$g(\mathbf{x}) = \frac{A}{K} \sum_{i=1}^K (\mathbf{w}_i) (\mathbf{x}; \mathbf{w}_i) \quad \epsilon;$$

Taking $\epsilon = 1/2$ (for which $K = O(A^2 B^2 \log N(\frac{1}{3AL}))$), by the probabilistic method, there exists a deterministic set $\{\mathbf{w}_i\}_{i \in [K]} \subset W$ and $f_i := (\mathbf{w}_i) g_{i \in [K]} \geq \int_W (\cdot; \mathbf{w}_i) d\mathbf{w}$ such that the above holds.

Next, note that both g (by (9)) and the function $\mathbf{x} \mapsto \frac{A}{K} \sum_{i=1}^K (\mathbf{w}_i) (\mathbf{x}; \mathbf{w}_i)$ are (AL) -Lipschitz. Therefore, for any $\mathbf{x} \in X$, taking \mathbf{x} to be the point in the covering set with $d_X(\mathbf{x}; \mathbf{x}) \leq \frac{\epsilon}{3AL}$, we have

$$\begin{aligned} g(\mathbf{x}) &= \frac{A}{K} \sum_{i=1}^K (\mathbf{w}_i) (\mathbf{x}; \mathbf{w}_i) \\ jg(\mathbf{x}) &= g(\mathbf{x}) + g(\mathbf{x}) = \frac{A}{K} \sum_{i=1}^K (\mathbf{w}_i) (\mathbf{x}; \mathbf{w}_i) + \frac{A}{K} \sum_{i=1}^K (\mathbf{w}_i) (\mathbf{x}; \mathbf{w}_i) = \frac{A}{K} \sum_{i=1}^K (\mathbf{w}_i) (\mathbf{x}; \mathbf{w}_i) \\ AL \frac{\epsilon}{3AL} &+ \frac{\epsilon}{3} + AL \frac{\epsilon}{3AL} = \epsilon. \end{aligned}$$

This proves the lemma. □

Proposition B.1 (Approximating smooth k -variable functions). *For any $\epsilon_{\text{approx}} > 0$, $R \geq 1$, $C > 0$, we have the following: Any $(R; C)$ -smooth function (Definition B.1) $g : \mathbb{R}^k \rightarrow \mathbb{R}$ is $(\epsilon_{\text{approx}}; R; M; C)$ -approximable by sum of relus (Definition D.1) with $M = C(k)C^2 \log(1 + C/\epsilon_{\text{approx}})$ and $C = C(k)C$, where $C(k) > 0$ is a constant that depends only on k . In other words, there exists*

$$f(\mathbf{z}) = \sum_{m=1}^M c_m (\mathbf{a}_m^\top [\mathbf{z}; 1]) \quad \text{with} \quad \sum_{m=1}^M |c_m| \leq C; \quad \max_{m \in [M]} |c_m| \leq 1;$$

such that $\sup_{\mathbf{z} \in [-R; R]^k} |f(\mathbf{z}) - g(\mathbf{z})| \leq \epsilon_{\text{approx}}$.

Proof. As function $g : \mathbb{B}_\infty^k(R) \rightarrow \mathbb{R}$ is $(R; C)$ -smooth, we can apply Lemma B.5 to obtain that there exists a signed measure ν over $W := \{\mathbf{w} \in \mathbb{R}^{k+1} : \|\mathbf{w}\|_1 = 1\}$ such that

$$g(\mathbf{z}) = \int_W \frac{1}{R} (\mathbf{w}^\top [\mathbf{z}; R]) d\nu(\mathbf{w}); \quad \forall \mathbf{z} \in [-R; R]^k;$$

and $A = \text{TV}(\nu) \leq C(k)C$ where $C(k) > 0$ denotes a constant depending only on k .

We now apply Lemma B.6 to approximate the above random feature by finitely many neurons. Let $\mathbf{x} := [\mathbf{z}; R] \in \mathbb{X} := [-R; R]^{k+1}$. Then, the function $\phi(\mathbf{x}; \mathbf{w}) := \frac{1}{R} (\mathbf{w}^\top \mathbf{x}) = \frac{1}{R} (\mathbf{w}^\top [\mathbf{z}; R])$ is bounded by $B = 1$ and $(1/R)$ -Lipschitz in \mathbf{x} (in the standard ℓ_∞ -distance). Further, we have $\log N(\mathbb{X}; k, k_\infty; \frac{\epsilon}{3A=R}) \leq O(k \log(1 + A/\epsilon))$. We can thus apply Lemma B.6 to obtain that, for

$$M = O(kA^2 \log(1 + A/\epsilon)) = \frac{2}{\epsilon^2} \log(1 + C/\epsilon) = \frac{2}{\epsilon^2} \log(1 + C/\epsilon);$$

there exists $\{g_m\}_{m \in [M]} \subset \mathcal{F}$ and $\mathbf{W} = \{\mathbf{w}_m\}_{m \in [M]} \subset W = \{\mathbf{w} \in \mathbb{R}^{k+1} : \|\mathbf{w}\|_1 = 1\}$ such that

$$\sup_{\mathbf{z} \in [-R; R]^k} |g(\mathbf{z}) - \frac{1}{M} \sum_{m=1}^M g_m(\mathbf{z})| \leq \frac{A}{M} \leq \frac{A}{M} \leq \epsilon;$$

where (recalling $\mathbf{z} = [S; \mathbf{f}]$)

$$g_m(\mathbf{z}) = \frac{A}{M} \sum_{m=1}^M \frac{1}{R} \mathbf{w}_m^\top [\mathbf{z}; R] = \frac{A}{M} \sum_{m=1}^M \frac{1}{R} \sum_{i=1}^k w_{m,i} z_i + \frac{A}{M} w_{m,k+1} = \frac{1}{M} \sum_{m=1}^M \left(\sum_{i=1}^k w_{m,i} z_i + w_{m,k+1} \right);$$

Note that we have $\sum_{m=1}^M \|\mathbf{w}_m\|_1 = A \leq C(k)C$, and $\sum_{m=1}^M \|\mathbf{w}_m\|_1 = 1$. This is the desired result. \square

B.3 Optimization

The following convergence result for minimizing a smooth and strongly convex function is standard from the convex optimization literature, see e.g. Bubeck [13, Theorem 3.10].

Proposition B.2 (Gradient descent for smooth and strongly convex functions). *Suppose $L : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly convex and L -smooth for some $0 < \mu \leq L$. Then, the gradient descent iterates $\mathbf{w}_{\text{GD}}^{t+1} := \mathbf{w}_{\text{GD}}^t - \eta \nabla L(\mathbf{w}_{\text{GD}}^t)$ with learning rate $\eta = 1/L$ and initialization $\mathbf{w}_{\text{GD}}^0 \in \mathbb{R}^d$ satisfies for any $t \geq 1$,*

$$\|\mathbf{w}_{\text{GD}}^t - \mathbf{w}^*\|_2 \leq \exp(-\mu t) \|\mathbf{w}_{\text{GD}}^0 - \mathbf{w}^*\|_2;$$

$$L(\mathbf{w}_{\text{GD}}^t) - L(\mathbf{w}^*) \leq \frac{\mu}{2} \exp(-\mu t) \|\mathbf{w}_{\text{GD}}^0 - \mathbf{w}^*\|_2^2;$$

where $\kappa := L/\mu$ is the condition number of L , and $\mathbf{w}^* := \arg \min_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w})$ is the minimizer of L .

The following convergence result of proximal gradient descent (PGD) on convex composite minimization problem is also standard, see e.g. [8].

Proposition B.3 (Proximal gradient descent for convex function). *Suppose $L = f + h$, $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex and L -smooth for some $L > 0$, $h : \mathbb{R}^d \rightarrow \mathbb{R}$ is a simple convex function. Then, the proximal gradient iterates $\mathbf{w}_{\text{PGD}}^{t+1} := \text{prox}_h(\mathbf{w}_{\text{PGD}}^t - \eta \nabla f(\mathbf{w}_{\text{PGD}}^t))$ with learning rate $\eta = 1/L$ and initialization $\mathbf{w}_{\text{GD}}^0 \in \mathbb{R}^d$ satisfies the following for any $t \geq 1$:*

- $fL(\mathbf{w}_{\text{PGD}}^t)$ is a decreasing sequence.
- For any minimizer $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^d} L(\mathbf{w})$,

$$L(\mathbf{w}_{\text{GD}}^{t+1}) - L(\mathbf{w}^*) \leq \frac{1}{2} \|\mathbf{w}_{\text{PGD}}^t - \mathbf{w}^*\|_2^2 + \|\mathbf{w}_{\text{PGD}}^{t+1} - \mathbf{w}^*\|_2^2;$$

and hence $\|\mathbf{w}_{\text{PGD}}^t - \mathbf{w}^*\|_2^2$ is also a decreasing sequence.

- For $k \geq 1; t \geq 0$, it holds that

$$L(\mathbf{w}_{\text{GD}}^{t+k}) - L(\mathbf{w}^*) \leq \frac{1}{2k} \|\mathbf{w}_{\text{PGD}}^t - \mathbf{w}^*\|_2^2;$$

B.4 Uniform convergence

The following result is shown in [87, Section 5.6].

Theorem B.1. Suppose that $\phi : [0; +\infty) \rightarrow [0; +\infty)$ is a convex, non-decreasing function that satisfies $\phi(x+y) \leq \phi(x) + \phi(y)$. For any random variable X , we consider the Orlicz norm induced by $\|\cdot\|_\phi := \inf \{K > 0 : \mathbb{E}(\phi(X/K)) \leq 1\}$.

Suppose that $\{f_X(t)\}$ is a zero-mean random process indexed by \mathcal{Z} such that $\|X\|_\phi \leq K$ and $\phi(\cdot)$ for some metric d on the space \mathcal{Z} . Then it holds that

$$\mathbb{P} \left(\sup_{j \in \mathcal{J}} |f_X(t_j) - f_X(t)| \geq \frac{1}{\phi(t/D)} \delta t \right) \leq \frac{1}{\phi(t/D)} \delta t$$

where D is the diameter of the metric space (\mathcal{Z}, d) , and the generalized Dudley entropy integral J is given by

$$J := \int_0^D \phi^{-1}(N(\cdot; \cdot; r)) dr$$

where $N(\cdot; \cdot; r)$ is the r -covering number of (\mathcal{Z}, d) .

As a corollary of Theorem B.1, we have the following result.

Proposition B.4 (Uniform concentration bound by chaining). Suppose that $\{f_X(t)\} \in \mathcal{B}_\phi$ is a zero-mean random process given by

$$X := \frac{1}{N} \sum_{i=1}^N f(z_i)$$

where z_1, \dots, z_N are i.i.d samples from a distribution \mathbb{P}_Z such that the following assumption holds:

- (a) The index set \mathcal{Z} is equipped with a distance d and diameter D . Further, assume that for some constant A , for any ball B_r of radius r in \mathcal{Z} , the covering number admits upper bound $\log N(\cdot; \cdot; r) \leq d \log(2Ar/r)$ for all $0 < r \leq D/2$.
- (b) For any fixed \mathcal{Z} and z sampled from \mathbb{P}_Z , the random variable $f(z)$ is a $\text{SG}(B^0)$ -sub-Gaussian random variable.
- (c) For any $\mathcal{Z}' \subseteq \mathcal{Z}$ and z sampled from \mathbb{P}_Z , the random variable $f(z) - f(z')$ is a $\text{SG}(B^1(\cdot; \cdot))$ -sub-Gaussian random variable.

Then with probability at least $1 - \epsilon$, it holds that

$$\sup_{j \in \mathcal{J}} |f_X(t_j) - f_X(t)| \leq C B^0 \sqrt{\frac{d \log(2A) + \log(1/\epsilon)}{N}}$$

where C is a universal constant, and we denote $B^0 = 1 + B^1 D = B^0$.

Furthermore, if we replace the SG in assumption (b) and (c) by SE, then with probability at least $1 - \epsilon$, it holds that

$$\sup_{j \in \mathcal{J}} |f_X(t_j) - f_X(t)| \leq C B^0 \sqrt{\frac{d \log(2A) + \log(1/\epsilon)}{N} + \frac{d \log(2A) + \log(1/\epsilon)}{N}}$$

Proof. Fix a $D_0 \in (0; D]$ to be specified later. We pick a $(D_0/2)$ -covering \mathcal{Z}_0 of \mathcal{Z} so that $\log |\mathcal{Z}_0| \leq d \log(2AD_0)$. Then, by the standard uniform covering of independent sub-Gaussian random variables, we have with probability at least $1 - \epsilon/2$,

$$\sup_{j \in \mathcal{Z}_0} |f_X(t_j) - f_X(t)| \leq C B^0 \sqrt{\frac{d \log(2AD_0) + \log(2/\epsilon)}{N}}$$

Assume that $\rho = \frac{1}{n}$; $\epsilon \in [0, \frac{1}{n}]$. For each $j \in [n]$, we consider B_j is the ball centered at x_j of radius D_0 in $(\mathbb{R}^D; \|\cdot\|)$. Then $\bigcup_{j \in [n]} B_j$ has diameter D_0 and admits covering number bound $\log N(\epsilon; \bigcup_{j \in [n]} B_j) \leq \frac{1}{\epsilon} \frac{D_0}{D}$. Hence, we can apply Theorem B.1 with the process $\sum_{j \in [n]} \langle X_j, \cdot \rangle$, then

$$\mathbb{P} \left(\left| \sum_{j \in [n]} \langle X_j, \cdot \rangle \right| \leq \frac{B^1}{N} \left(\frac{1}{\epsilon}; \cdot \right) \right);$$

and a simple calculation yields

$$\mathbb{P} \left(\sup_{i \in [n]} \left| \sum_{j \in [n]} \langle X_j, X_i \rangle \right| \leq C' B^1 D_0 \left(\frac{1}{\epsilon} \frac{D_0}{D} + t \right) \leq 2 \exp(-N t^2) \leq \delta \right);$$

Therefore, we can let $t = \frac{\sqrt{\log(2n/\delta)}}{\sqrt{N}}$ in the above inequality and taking the union bound over $j \in [n]$, and hence with probability at least $1 - \delta$, it holds that for all $j \in [n]$,

$$\sup_{i \in [n]} \left| \sum_{j \in [n]} \langle X_j, X_i \rangle \right| \leq C' B^1 D_0 \frac{\sqrt{2d \log(2AD=D_0) + \log(4/\delta)}}{N};$$

Notice that for each $i \in [n]$, there exists $j \in [n]$ such that $i \in B_j$, and hence

$$\left| \sum_{j \in [n]} \langle X_j, X_i \rangle \right| \leq \sum_{j \in [n]} \langle X_j, X_i \rangle;$$

Thus, with probability at least $1 - \delta$, it holds

$$\sup_{i \in [n]} \left| \sum_{j \in [n]} \langle X_j, X_i \rangle \right| \leq \sup_{i \in [n]} \sum_{j \in [n]} \langle X_j, X_i \rangle \leq C'' (B_0 + B^1 D_0) \frac{\sqrt{d \log(2AD=D_0) + \log(2/\delta)}}{N};$$

Taking $D_0 = D$ completes the proof of SG case.

We next consider the SE case. The idea is the same as the SG case, but in this case we need to consider the following Orlicz-norm:

$$\psi(t) := \exp \left(\frac{N t^2}{t+1} \right) - 1;$$

Then Bernstein's inequality of SE random variables yields

$$\mathbb{P} \left(\left| \sum_{j \in [n]} \langle X_j, X_i \rangle \right| \leq C_0 B^1 \left(\frac{1}{\epsilon}; \cdot \right) \right);$$

for some universal constant C_0 . Therefore, we can repeat the argument above to deduce that with probability at least $1 - \delta$, it holds

$$\sup_{i \in [n]} \left| \sum_{j \in [n]} \langle X_j, X_i \rangle \right| \leq C'' (B_0 + B^1 D_0) \frac{\sqrt{d \log(2AD=D_0) + \log(2/\delta)}}{N} + \frac{d \log(2AD=D_0) + \log(2/\delta)}{N};$$

Taking $D_0 = D$ completes the proof. \square

B.5 Useful properties of transformers

The following result can be obtained immediately by “joining” the attention heads and MLP layers of two single-layer transformers.

Proposition B.5 (Joining parallel single-layer transformers). *Suppose that $P_1 : \mathbb{R}^{(D_0+D_1) \times N} \rightarrow \mathbb{R}^{D_1 \times N}$; $P_2 : \mathbb{R}^{(D_0+D_2) \times N} \rightarrow \mathbb{R}^{D_2 \times N}$ are two sequence-to-sequence functions that are implemented by single-layer transformers, i.e. there exists $\mathcal{H}_1, \mathcal{H}_2$ such that*

$$\begin{aligned} \text{TF}_1 : \mathbf{H}_1 &= \begin{matrix} \mathbf{h}_i^{(0)} \\ \mathbf{h}_i^{(1)} \end{matrix}_{1 \leq i \leq N} \xrightarrow{P_1} \mathbf{H}_1^{(0)}; \\ \text{TF}_2 : \mathbf{H}_2 &= \begin{matrix} \mathbf{h}_i^{(0)} \\ \mathbf{h}_i^{(2)} \end{matrix}_{1 \leq i \leq N} \xrightarrow{P_2} \mathbf{H}_2^{(0)}; \end{aligned}$$

Then, there exists \mathbf{H}' such that for \mathbf{H}' that takes form $\mathbf{h}'_i = [\mathbf{h}_i^{(0)}; \mathbf{h}_i^{(1)}; \mathbf{h}_i^{(2)}]$, with $\mathbf{h}_i^{(0)} \in \mathbb{R}^{D_0}; \mathbf{h}_i^{(1)} \in \mathbb{R}^{D_1}; \mathbf{h}_i^{(2)} \in \mathbb{R}^{D_2}$, we have

$$\text{TF} : \mathbf{H}' = \begin{matrix} \begin{matrix} \mathbf{h}_i^{(0)} \\ \mathbf{h}_i^{(1)} \\ \mathbf{h}_i^{(2)} \end{matrix} \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ \begin{matrix} 2 \\ 4 \\ 1 \end{matrix} \end{matrix} \begin{matrix} \begin{matrix} \mathbf{H}^{(0)} \\ \mathbf{H}_1 \\ \mathbf{H}_2 \end{matrix} \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ \begin{matrix} 2 \\ 4 \\ 2 \end{matrix} \end{matrix} \in \mathbb{R}^{(D_0+D_1+D_2) \times N} \xrightarrow{\text{---}} \begin{matrix} \mathbf{H}^{(0)} \\ P_1(\mathbf{H}_1) \\ P_2(\mathbf{H}_2) \end{matrix} \quad \#$$

Further, \mathbf{H}' has at most $M = M_1 + M_2$ heads, $D' = D_1 + D_2$ hidden dimension in its MLP layer, and norm bound $\|\mathbf{H}'\| \leq \|\mathbf{H}_1\| + \|\mathbf{H}_2\|$.

Proposition B.6 (Joining parallel multi-layer transformers). *Suppose that $P_1 : \mathbb{R}^{(D_0+D_1) \times N} \rightarrow \mathbb{R}^{D_1 \times N}; P_2 : \mathbb{R}^{(D_0+D_2) \times N} \rightarrow \mathbb{R}^{D_2 \times N}$ are two sequence-to-sequence functions that are implemented by multi-layer transformers, i.e. there exists $\mathbf{H}_1; \mathbf{H}_2$ such that*

$$\begin{aligned} \text{TF}_1 : \mathbf{H}_1 &= \begin{matrix} \mathbf{h}_i^{(0)} \\ \mathbf{h}_i^{(1)} \\ \mathbf{h}_i^{(2)} \end{matrix} \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \begin{matrix} \mathbf{H}^{(0)} \\ P_1(\mathbf{H}_1) \\ P_2(\mathbf{H}_2) \end{matrix} \in \mathbb{R}^{(D_0+D_1) \times N} \xrightarrow{\text{---}} \begin{matrix} \mathbf{H}^{(0)} \\ P_1(\mathbf{H}_1) \\ P_2(\mathbf{H}_2) \end{matrix} \quad \# \\ \text{TF}_2 : \mathbf{H}_2 &= \begin{matrix} \mathbf{h}_i^{(0)} \\ \mathbf{h}_i^{(1)} \\ \mathbf{h}_i^{(2)} \end{matrix} \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \begin{matrix} \mathbf{H}^{(0)} \\ P_1(\mathbf{H}_1) \\ P_2(\mathbf{H}_2) \end{matrix} \in \mathbb{R}^{(D_0+D_2) \times N} \xrightarrow{\text{---}} \begin{matrix} \mathbf{H}^{(0)} \\ P_1(\mathbf{H}_1) \\ P_2(\mathbf{H}_2) \end{matrix} \quad \# \end{aligned}$$

Then, there exists \mathbf{H}' such that for \mathbf{H}' that takes form $\mathbf{h}'_i = [\mathbf{h}_i^{(0)}; \mathbf{h}_i^{(1)}; \mathbf{h}_i^{(2)}]$, with $\mathbf{h}_i^{(0)} \in \mathbb{R}^{D_0}; \mathbf{h}_i^{(1)} \in \mathbb{R}^{D_1}; \mathbf{h}_i^{(2)} \in \mathbb{R}^{D_2}$, we have

$$\text{TF} : \mathbf{H}' = \begin{matrix} \begin{matrix} \mathbf{h}_i^{(0)} \\ \mathbf{h}_i^{(1)} \\ \mathbf{h}_i^{(2)} \end{matrix} \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ \begin{matrix} 2 \\ 4 \\ 1 \end{matrix} \end{matrix} \begin{matrix} \begin{matrix} \mathbf{H}^{(0)} \\ \mathbf{H}_1 \\ \mathbf{H}_2 \end{matrix} \\ \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ \begin{matrix} 2 \\ 4 \\ 2 \end{matrix} \end{matrix} \in \mathbb{R}^{(D_0+D_1+D_2) \times N} \xrightarrow{\text{---}} \begin{matrix} \mathbf{H}^{(0)} \\ P_1(\mathbf{H}_1) \\ P_2(\mathbf{H}_2) \end{matrix} \quad \#$$

Further, \mathbf{H}' has at most $L = \max\{L_1; L_2\}$ layers, $\max_{i \in [L]} M^{(i)} = \max_{i \in [L]} M_1^{(i)} + M_2^{(i)}$ heads, $\max_{i \in [L]} D^{(i)} = \max_{i \in [L]} D_1^{(i)} + D_2^{(i)}$ hidden dimension in its MLP layer (understanding the size of the empty layers as 0), and norm bound $\|\mathbf{H}'\| \leq \|\mathbf{H}_1\| + \|\mathbf{H}_2\|$.

Proof. When $L_1 = L_2$ (\mathbf{H}_1 and \mathbf{H}_2 have the same number of layers), the result follows directly by applying Proposition B.5 repeatedly for all L_1 layers and the definition of the norm (2).

If (without loss of generality) $L_1 < L_2$, we can augment \mathbf{H}_1 to L_2 layers by adding $(L_2 - L_1)$ layers with zero attention heads, and zero MLP hidden dimension (note that this does not change M_1, D_1 , and $\|\mathbf{H}_1\|$). Due to the residual structure, the transformer maintains the output $P_1(\mathbf{H}_1)$ throughout layer $L_1 + 1; \dots; L_2$, and it reduces to the case $L_1 = L_2$. \square

C Extension to decoder-based architecture

Here we briefly discuss how our theoretical results can be adapted to decoder-based architectures (henceforth decoder TFs). Adopting the setting as in Section 2, we consider a sequence of N input vectors $\{\mathbf{h}_i\}_{i=1}^N \in \mathbb{R}^D$, written compactly as an input matrix $\mathbf{H} = [\mathbf{h}_1; \dots; \mathbf{h}_N] \in \mathbb{R}^{D \times N}$. Recall that $\text{ft}(t) := \text{ReLU}(t) = \max\{t; 0\}$ denotes the standard relu activation.

C.1 Decoder-based transformers

Decoder TFs are the same as encoder TFs, except that the attention layers are replaced by masked attention layers with a specific decoder-based (causal) attention mask.

Definition C.1 (Masked attention layer). *A masked attention layer with M heads is denoted as $\text{MAttn}(\cdot)$ with parameters $\{\mathbf{V}_m; \mathbf{Q}_m; \mathbf{K}_m\}_{m \in [M]} \in \mathbb{R}^{D \times D}$. On any input sequence*

$$\mathbf{H} \in \mathbb{R}^{D \times N^0} \text{ with } N^0 \leq N, \quad \mathbf{H} = \text{MAttn}(\mathbf{H}) := \mathbf{H} + \prod_{m=1}^M (\mathbf{V}_m \mathbf{H}) \text{ (MSK}_{1:N^0; 1:N^0}) (\mathbf{Q}_m \mathbf{H})^\top (\mathbf{K}_m \mathbf{H}) \in \mathbb{R}^{D \times N^0}; \quad (10)$$

where \odot denotes the entry-wise (Hadamard) product of two matrices, and $\text{MSK} \in \mathbb{R}^{N \times N}$ is the mask matrix given by

$$\text{MSK} = \begin{bmatrix} 1 & 1=2 & 1=3 & \dots & 1=N \\ 0 & 1=2 & 1=3 & \dots & 1=N \\ 0 & 0 & 1=3 & \dots & 1=N \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1=N \end{bmatrix}.$$

In vector form, we have

$$\hat{\mathbf{h}}_i = [\text{Attn}(\mathbf{H})]_i = \mathbf{h}_i + \prod_{m=1}^M \frac{1}{i} \prod_{j=1}^i (g_m \mathbf{Q}_m \mathbf{h}_i; \mathbf{K}_m \mathbf{h}_j) \mathbf{V}_m \mathbf{h}_j.$$

Notice that standard masked attention definitions use the pre-activation additive masks (with mask value -1) [84]. The post-activation multiplicative masks we use is equivalent to the pre-activation additive masks, and the modified presentation is for notational convenience. We also use a normalized ReLU activation $f(t) = i$ in place of the standard softmax activation to be consistent with Definition 1. Note that the normalization $1=i$ is to ensure that the attention weights $f(g_m \mathbf{Q}_m \mathbf{h}_i; \mathbf{K}_m \mathbf{h}_j) = i g_{j \in [i]}$ is a set of non-negative weights that sum to $O(1)$. The motivation of masked attention layer is to ensure that, when processing a sequence of tokens, the computations at any token do not see any later token.

We next define the decoder-based transformers with $L - 1$ transformer layers, each consisting of a masked attention layer (c.f. Definition C.1) followed by an MLP layer (c.f. Definition 2). This definition is similar to the definition of encoder-based transformers (c.f., Definition 3), except that we replace the attention layers by masked attention layers.

Definition C.2 (Decoder-based Transformer). *An L -layer decoder-based transformer, denoted as $\text{DTF}(\cdot)$, is a composition of L self-attention layers each followed by an MLP layer: $\mathbf{H}^{(L)} = \text{DTF}(\mathbf{H}^{(0)})$, where $\mathbf{H}^{(0)} \in \mathbb{R}^{D \times N}$ is the input sequence, and*

$$\mathbf{H}^{(l)} = \text{MLP}_{\text{mlp}}^{(l)} \circ \text{MAttn}_{\text{mattn}}^{(l)} \circ \mathbf{H}^{(l-1)}; \quad l \in \{1, \dots, L\}.$$

Above, the parameter $\theta^{(l)} = (\theta_{\text{mattn}}^{(l)}, \theta_{\text{mlp}}^{(l)})$ is the parameter consisting of the attention layers $\theta_{\text{mattn}}^{(l)} = f(\mathbf{V}_m^{(l)}; \mathbf{Q}_m^{(l)}; \mathbf{K}_m^{(l)}) g_{m \in [M^{(l)}]} \in \mathbb{R}^{D \times D}$ and the MLP layers $\theta_{\text{mlp}}^{(l)} = (\mathbf{W}_1^{(l)}; \mathbf{W}_2^{(l)}) \in \mathbb{R}^{D^{(l)} \times D} \times \mathbb{R}^{D \times D^{(l)}}$. We will frequently consider ‘‘attention-only’’ decoder-based transformers with $\mathbf{W}_1^{(l)}; \mathbf{W}_2^{(l)} = \mathbf{0}$, which we denote as $\text{DTF}^0(\cdot)$ for shorthand, with $\theta^{(l)} = (1:L) := \theta_{\text{mattn}}^{(l)}$.

We also use (2) to define the norm of $\text{DTF}(\cdot)$.

C.2 In-context learning with decoder-based transformers

We consider using decoder-based TFs to perform ICL. We encode $(D; \mathbf{x}_{N+1})$, which follows the generating rule as described in Section 2.2, into an input sequence $\mathbf{H} \in \mathbb{R}^{D \times (2N+1)}$. In our theory, we use the following format, where the first two rows contain $(D; \mathbf{x}_{N+1})$ which alternating between $[\mathbf{x}_i; 0] \in \mathbb{R}^{d+1}$ and $[0_{d \times 1}; y_i] \in \mathbb{R}^{d+1}$ (the same setup as adopted in [31, 2]); The third row contains fixed vectors $\{\mathbf{p}_i\}_{i \in [N+1]}$ with ones, zeros, the example index, and indicator for being the covariate token (similar to a positional encoding vector):

$$\mathbf{H} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{0} & \dots & \mathbf{x}_N & \mathbf{0} & \mathbf{x}_{N+1} \\ 0 & y_1 & \dots & 0 & y_N & 0 \\ \mathbf{p}_1 & \mathbf{p}_2 & \dots & \mathbf{p}_{2N-1} & \mathbf{p}_{2N} & \mathbf{p}_{2N+1} \end{bmatrix}; \quad \mathbf{p}_i := \begin{cases} 0_{D-(d+4)} & \text{if } i \equiv 2 \pmod{2} \\ 1 & \text{if } i \equiv 1 \pmod{2} \\ i & \text{if } i \equiv N+1 \pmod{2} \\ \text{mod}(i+1; 2) & \text{otherwise} \end{cases} \in \mathbb{R}^{D-(d+1)}; \quad (11)$$

(11) is different from our input format (3) for encoder-based TFs. The main difference is that $(\mathbf{x}_i; y_i)$ are in different tokens in (11), whereas $(\mathbf{x}_i; y_i)$ are in the same token in (3). The reason for the former (i.e., different tokens in decoder) is that we want to avoid every $[\mathbf{x}_i; 0]$ token seeing the information of y_i , since we will evaluate the loss at every token. The reason for the latter (i.e., the same token in encoder) is for presentation convenience: since we only evaluate the loss at the last token, it is not necessary to alternate between $[\mathbf{x}_i; 0]$ and $[0; y_i]$ to avoid information leakage.

We then feed \mathbf{H} into a decoder TF to obtain the output $\hat{\mathbf{H}} = \text{DTF}(\mathbf{H}) \in \mathbb{R}^{D \times (2N+1)}$ with the same shape, and *read out* the prediction \hat{y}_{N+1} from the $(d+1; 2N+1)$ -th entry of $\hat{\mathbf{H}} = [\hat{\mathbf{h}}_i]_{i \in [2N+1]}$ (the entry corresponding to the last missing test label): $\hat{y}_{N+1} = \text{read}_y(\hat{\mathbf{H}}) := (\hat{\mathbf{h}}_{2N+1})_{d+1}$. The goal is to predict \hat{y}_{N+1} that is close to $y_{N+1} \sim \mathcal{P}_{y|x_{N+1}}$ measured by proper losses.

The benefit of using the decoder architecture is that, during the pre-training phase, one can construct the training loss function by using all the predictions $\hat{y}_j, j \in [N+1]$, where \hat{y}_j gives the $(d+1; 2j-1)$ -th entry of $\hat{\mathbf{H}} = [\hat{\mathbf{h}}_i]_{i \in [2N+1]}$ for each $j \in [N+1]$ (the entry corresponding to the missing test label of the $2j-1$ 'th token): $\hat{y}_j = \text{read}_{y,j}(\hat{\mathbf{H}}) := (\hat{\mathbf{h}}_{2j-1})_{d+1}$. Given a loss function $\ell: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ associated to a single response, the training loss associated to the whole input sequence can be defined by $\ell(\mathbf{H}) = \sum_{j=1}^{N+1} \ell(y_j; \hat{y}_j)$. This potentially enables less training sequences in the pre-training stage, and some generalization bound analysis justifying this benefit was provided in [47].

C.3 Results

We discuss how our theoretical results upon encoder TFs can be converted to those of the decoder TFs. Taking the implementation of (ICGD) (a key mechanism that enables most basic ICL algorithms such as ridge regression; cf. Appendix D.1) as an example, this conversion is enabled by the following facts: (a) the input format (11) of decoders can be converted to the input format (3) of encoders by a 2-layer decoder TF; (b) the encoder TF that implements (ICGD) with input format (3), by a slight parameter modification, can be converted to a decoder TF that implements the (ICGD) algorithm with a converted input format.

Input format conversion Despite the difference between the input format (11) and (3), we show that there exists a 2-layer decoder TF that can convert the input format (11) to format (3). The proof can be found in Appendix C.4.

Proposition C.1 (Input format conversion). *There exists a 2-layer decoder TF DTF with 3 heads per layer, hidden dimension 2 and $\# \# = 12$ such that upon taking input \mathbf{H} of format (11), it outputs $\hat{\mathbf{H}} = \text{DTF}(\mathbf{H})$ with*

$$\hat{\mathbf{H}} = \begin{matrix} & \mathbf{x}_1 & \mathbf{x}_1 & \dots & \mathbf{x}_N & \mathbf{x}_N & \mathbf{x}_{N+1} & \# \\ \mathbf{H} = & 0 & y_1 & \dots & 0 & y_N & 0 & : \\ & \mathbf{p}_1 & \mathbf{p}_2 & \dots & \mathbf{p}_{2N-1} & \mathbf{p}_{2N} & \mathbf{p}_{2N+1} & \end{matrix} \quad (12)$$

In particular, format (12) contains format (3) as a submatrix, by restricting to the $f1; 2; \dots; D-1; D-2; Dg$ rows and $f2; 4; \dots; 2N-2; 2N; 2N+1g$ columns.

Generalization TF constructions to decoder architecture The construction in Theorem D.1 can be generalized to using the input format (12) along with a decoder TF, by using the scratch pad within the last token to record the gradient descent iterates. Further, if we slightly change the normalization in MSK from $1=i$ to $1=((i-1)-1)$, then the same construction performs (ICGD) (with training examples $f1; \dots; jg$) at every token $i = 2j+1$ (corresponding to predicting at \mathbf{x}_{j+1}). Building on this extension, all our constructions in Section 3 and Section 4.2 can be generalized to decoder TFs.

C.4 Proof of Proposition C.1

For the simplicity of presentation, we write $c_i = di=2e; t_i = \text{mod}(i+1; 2)$, $\mathbf{u}_i = \mathbf{h}_i[1:d] \in \mathbb{R}^{d+1}$ be the vector of first d entries of \mathbf{h}_i ⁶, and let $v_i = \mathbf{h}_i[d+1]$ be the $(d+1)$ -th entry of \mathbf{h}_i . With such notations, the input sequence $\mathbf{H} = [\mathbf{h}_i]_i$ can be compactly written as

$$\mathbf{h}_i = [\mathbf{u}_i; v_i; \mathbf{0}_{D-d-4}; c_i; 1; t_i]:$$

In the following, we construct the desired $\hat{\mathbf{H}} = (1); (2)$ as follows.

Step 1: construction of $(1) = (1)_{\text{mattn}}; (1)_{\text{mlp}}$, so that $\text{MLP}_{\text{mlp}}^{(1)} \text{Mattn}_{\text{mattn}}^{(1)}$ maps

$$\mathbf{h}_i \xrightarrow{\text{Mattn}_{\text{mattn}}^{(1)}} \mathbf{h}'_i = [\mathbf{u}_i; v_i; \mathbf{0}_{D-d-6}; t_i(c_i^2 + 0.5); t_i c_i; c_i; 1; t_i]$$

⁶In other words, when $2-i$, $\mathbf{u}_i = \mathbf{x}_{(i-1)-2}$; when $2 \nmid i$, $\mathbf{u}_i = \mathbf{0}_d$.

$$\text{MLP}_{\text{mlp}}^{(1)} \mathbf{h}_i^{(1)} = [\mathbf{u}_i; v_i; \mathbf{0}_{D-d-6}; t_i c_i^2; t_i c_i; c_i; 1; t_i]:$$

For $m \geq 0; 1g$, we define matrices $\mathbf{Q}_m^{(1)}; \mathbf{K}_m^{(1)}; \mathbf{V}_m^{(1)} \in \mathbb{R}^{D \times D}$ such that

$$\mathbf{Q}_0^{(1)} \mathbf{h}_i = \mathbf{Q}_1^{(1)} \mathbf{h}_i = \begin{bmatrix} t_i \\ \mathbf{0} \end{bmatrix}; \quad \mathbf{K}_0^{(1)} \mathbf{h}_j = \mathbf{K}_1^{(1)} \mathbf{h}_j = \begin{bmatrix} c_j \\ \mathbf{0} \end{bmatrix}; \quad \mathbf{V}_0^{(1)} \mathbf{h}_j = \begin{bmatrix} \mathbf{0}_{D-4} \\ 3c_j \\ \mathbf{0}_3 \end{bmatrix}; \quad \mathbf{V}_1^{(1)} \mathbf{h}_j = \begin{bmatrix} \mathbf{0}_{D-3} \\ 2 \\ \mathbf{0}_2 \end{bmatrix};$$

for all $i; j$. By the structure of \mathbf{h}_i , these matrices indeed exist, and further it is straightforward to check that they have norm bounds

$$\max_m \|\mathbf{Q}_m^{(1)}\|_{\text{op}} \leq 1; \quad \max_m \|\mathbf{K}_m^{(1)}\|_{\text{op}} \leq 1; \quad \prod_m \|\mathbf{V}_m^{(1)}\|_{\text{op}} \leq 5:$$

Now, for every i ,

$$\frac{1}{i} \prod_{j=1}^i \sum_{m \in \{0,1\}} \mathbf{Q}_m^{(1)} \mathbf{h}_i; \mathbf{K}_m^{(1)} \mathbf{h}_j \quad \mathbf{V}_m^{(1)} \mathbf{h}_j = \frac{1}{i} \prod_{j=1}^i t_i [\mathbf{0}_{D-4}; 3c_j^2; 2c_j; \mathbf{0}; \mathbf{0}]:$$

Notice that $t_i \neq 0$ only when $2 \nmid j$, we then compute for $i = 2k$ that

$$\prod_{j=1}^i 3c_j^2 = 3 \frac{k(k-1)(2k-1)}{3} + 3k^2 = 2k^3 + k; \quad \prod_{j=1}^i 2c_j = 2 \frac{k(k-1)}{2} + 2k = 2k^2:$$

Therefore, the $\text{mattn}^{(1)} = f(\mathbf{Q}_m^{(1)}; \mathbf{K}_m^{(1)}; \mathbf{V}_m^{(1)})_{m \in \{0,1\}}$ we construct above is indeed the desired attention layer. The existence of the desired $\text{mlp}^{(1)}$ is clear, and $\text{mlp}^{(1)} = (\mathbf{W}_1^{(1)}; \mathbf{W}_2^{(1)})$ can further be chosen so that $\|\mathbf{W}_1^{(1)}\|_{\text{op}} \leq 1; \|\mathbf{W}_2^{(1)}\|_{\text{op}} \leq 1$.

Step 2: construction of $\text{mattn}^{(2)}$. For every $m \geq 0; 1; 0; 1g$, we define matrices $\mathbf{Q}_m^{(2)}; \mathbf{K}_m^{(2)}; \mathbf{V}_m^{(2)} \in \mathbb{R}^{D \times D}$ such that

$$\begin{aligned} \mathbf{Q}_0^{(2)} \mathbf{h}_i^{(1)} &= \mathbf{Q}_1^{(2)} \mathbf{h}_i^{(1)} = \mathbf{Q}_{-1}^{(2)} \mathbf{h}_i^{(1)} = \begin{bmatrix} 2t_i c_i^2 \\ t_i c_i \\ \mathbf{0} \end{bmatrix}; \\ \mathbf{K}_0^{(2)} \mathbf{h}_j^{(1)} &= \begin{bmatrix} c_j \\ \mathbf{0} \end{bmatrix}; \quad \mathbf{K}_1^{(2)} \mathbf{h}_j^{(1)} = \begin{bmatrix} c_j + 1 \\ \mathbf{0} \end{bmatrix}; \quad \mathbf{K}_{-1}^{(2)} \mathbf{h}_j^{(1)} = \begin{bmatrix} c_j - 1 \\ \mathbf{0} \end{bmatrix}; \\ \mathbf{V}_0^{(2)} \mathbf{h}_j^{(1)} &= \begin{bmatrix} 4\mathbf{u}_j \\ \mathbf{0}_{D-d} \end{bmatrix}; \quad \mathbf{V}_1^{(2)} \mathbf{h}_j^{(1)} = \mathbf{V}_{-1}^{(2)} \mathbf{h}_j^{(1)} = \begin{bmatrix} 2\mathbf{u}_j \\ \mathbf{0}_{D-d} \end{bmatrix}; \end{aligned}$$

for all $i; j$. By the structure of $\mathbf{h}_i^{(1)}$, these matrices indeed exist, and further it is straightforward to check that they have norm bounds

$$\max_m \|\mathbf{Q}_m^{(2)}\|_{\text{op}} \leq 1; \quad \max_m \|\mathbf{K}_m^{(2)}\|_{\text{op}} \leq 2; \quad \prod_m \|\mathbf{V}_m^{(2)}\|_{\text{op}} \leq 8:$$

Now, for every $i; j$, we have

$$\begin{aligned} & \prod_{m \in \{-1,0,1\}} \sum_{m \in \{-1,0,1\}} \mathbf{Q}_m^{(2)} \mathbf{h}_i^{(1)}; \mathbf{K}_m^{(2)} \mathbf{h}_j^{(1)} \quad \mathbf{V}_m^{(2)} \mathbf{h}_j^{(1)} \\ &= 2 \frac{t_i c_i^2}{c_i} \frac{t_i c_i}{c_j} + \frac{t_i c_i^2}{c_i} \frac{t_i c_i}{c_j + 1} + \frac{t_i c_i^2}{c_i} \frac{t_i c_i}{c_j - 1} \quad 2[\mathbf{u}_j; \mathbf{0}_{D-d}] \\ &= f \frac{2}{c_i} (c_i - c_j) + ((c_i - c_j) - 1) + ((c_i - c_j) + 1)g \quad 2c_i t_i [\mathbf{u}_j; \mathbf{0}_{D-d}] \\ &= (c_i - c_j) \quad 2c_i t_i [\mathbf{u}_j; \mathbf{0}_{D-d}]; \end{aligned}$$

where the last equality follows from the fact that

$$2 \frac{x}{x-1} + \frac{x}{x+1} + \frac{x}{x-1} = \begin{cases} < 0; & x = 1 \text{ or } x = -1; \\ x+1; & x \geq [1; 0]; \\ 1-x; & x \geq [0; 1]; \end{cases}$$

Therefore,

$$\begin{aligned} \prod_{j=1}^D \prod_{m \in \{-1, 0, 1\}} \mathbf{Q}_m^{(2)} \mathbf{h}_j^{(1)}; \mathbf{K}_m^{(2)} \mathbf{h}_j^{(1)} \mathbf{V}_m^{(2)} \mathbf{h}_j^{(1)} &= \prod_{j=1}^D 2 \mathbb{1}(c_i = c_j) c_i t_j [\mathbf{u}_j; \mathbf{0}_{D-d}] \\ &= \begin{cases} [\mathbf{x}_k; \mathbf{0}_{D-d}] & i = 2k \\ \mathbf{0}_D & \text{otherwise} \end{cases} \end{aligned}$$

Therefore, the $\text{mattn}^{(2)} = f(\mathbf{Q}_m^{(2)}; \mathbf{K}_m^{(2)}; \mathbf{V}_m^{(2)} \geq \mathbb{R}^{D \times D}) g_{m \in \{-1, 0, 1\}}$ we construct above maps

$$\mathbf{h}_i^{(1)} \mapsto \mathbf{h}_i'' = [\mathbf{x}_{\lceil i/2 \rceil}; v_i; \mathbf{0}_{D-d-6}; t_i c_i^2; t_i c_i; c_i; 1; t_i];$$

Finally, we only need to take a MLP layer $\text{mlp}^{(2)} = (\mathbf{W}_1^{(2)}; \mathbf{W}_2^{(2)})$ with hidden dimension 2 that maps

$$\mathbf{h}_i'' \mapsto \mathbf{h}_i^{(2)} = [\mathbf{x}_{\lceil i/2 \rceil}; v_i; \mathbf{0}_{D-d-6}; 0; 0; c_i; 1; t_i];$$

which clearly exists and can be chosen so that $k\mathbf{W}_1^{(2)} k_{\text{op}} \leq 1; k\mathbf{W}_2^{(2)} k_{\text{op}} \leq 1$.

Combining the two steps above, we complete the proof of Proposition C.1. \square

D Mechanism: In-context gradient descent

Technically, the constructions in Section 3.1-3.2 rely on a new efficient construction for transformers to implement in-context gradient descent and its variants, which we present as follows. We begin by presenting the result for implementing (vanilla) gradient descent on convex empirical risks.

Compact notation of input We will often use shorthand $y_i \geq \mathbb{R}$ defined as $y_i' = y_i$ for $i \geq [N]$ and $y_{N+1}' = 0$ to simplify our notation, with which the input sequence $\mathbf{H} \geq \mathbb{R}^{D \times (N+1)}$ can be compactly written as $\mathbf{h}_i = [\mathbf{x}_i; y_i'; \mathbf{p}_i] = [\mathbf{x}_i; y_i'; \mathbf{0}_{D-d-3}; 1; t_i]$ for $i \geq [N+1]$, where $t_i := 1/i < N+1/g$ is the indicator for the training examples.

D.1 Gradient descent on convex empirical risk

Let $\ell(\cdot; \cdot) : \mathbb{R}^2 \rightarrow \mathbb{R}$ be a loss function. Let $\mathbb{E}_N(\mathbf{w}) := \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{w}^\top \mathbf{x}_i; y_i)$ denote the empirical risk with loss function ℓ on dataset $f(\mathbf{x}_i; y_i) g_{i \in [N]}$, and

$$\mathbf{w}_{\text{GD}}^{t+1} := \mathbf{w}_{\text{GD}}^t - \eta \mathbb{E}_N(\mathbf{w}_{\text{GD}}^t) \quad (\text{ICGD})$$

denote the gradient descent trajectory on \mathbb{E}_N with initialization $\mathbf{w}_{\text{GD}}^0 \geq \mathbb{R}^d$ and learning rate $\eta > 0$.

We require the partial derivative of the loss $\partial_s \ell : (s; t) \mapsto \partial_s \ell(s; t)$ (as a bivariate function) to be approximable by a sum of relus, defined as follows.

Definition D.1 (Approximability by sum of relus). *A function $g : \mathbb{R}^k \rightarrow \mathbb{R}$ is (" $\text{approx}; R; M; C$ ")-approximable by sum of relus, if there exists a (" $M; C$ ")-sum of relus" function*

$$f_{M;C}(\mathbf{z}) = \sum_{m=1}^M c_m (\mathbf{a}_m^\top [\mathbf{z}; 1]) \quad \text{with} \quad \sum_{m=1}^M |c_m| \leq C; \max_{m \in [M]} \|\mathbf{a}_m\|_1 \leq 1; \mathbf{a}_m \geq \mathbb{R}^{k+1}; c_m \geq \mathbb{R};$$

such that $\sup_{\mathbf{z} \in [-R; R]^k} |g(\mathbf{z}) - f_{M;C}(\mathbf{z})| \leq \text{approx}$.

Definition D.1 is known to contain broad class of functions. For example, any mildly smooth k -variate function is approximable by a sum of relus for any (" $\text{approx}; R$ "), with mild bounds on $(M; C)$ (Proposition B.1, building on results of Bach [4]). Also, any function that is a $(M; C)$ -sum of relus itself (which includes all piecewise linear functions) is by definition $(0; 1; M; C)$ -approximable by sum of relus.

We show that L steps of (ICGD) can be approximately implemented by an $(L+1)$ -layer transformer.

Theorem D.1 (Convex ICGD). *Fix any $B_w > 0, L > 1, \eta > 0$, and $\text{approx} \leq B_w = (2L)$. Suppose that*

1. *The loss $\ell(\cdot; \cdot)$ is convex in the first argument;*
2. *$\partial_s \ell$ is (" $\text{approx}; R; M; C$ ")-approximable by sum of relus with $R = \max\{\eta B_x B_w; B_y; 1/g\}$.*

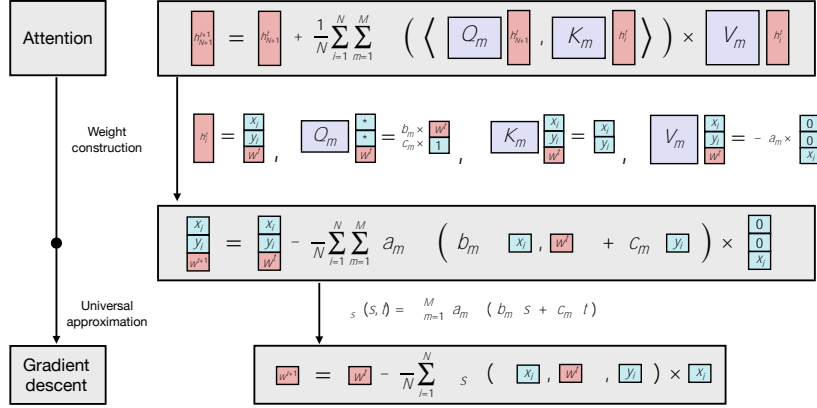


Figure 4: Illustration of our main mechanism for implementing basic ICL algorithms: One attention layer implements a single (ICGD) iterate (Proposition E.1 & Theorem D.1). Top: the attention mechanism as in Definition 1. Bottom: A single (ICGD) iterate. Middle: Linear algebraic illustration of the attention layer for implementing a GD update.

Then, there exists an attention-only transformer TF^0 with $(L + 1)$ layers, $\max_{\cdot \in [L]} M^{(\cdot)} = M$ heads within the first L layers, and $M^{(L+1)} = 2$ such that for any input data $(D; \mathbf{x}_{N+1})$ such that

$$\sup_{\|\mathbf{w}\|_2 \leq B_w} \max(r^{-2} \mathcal{L}_N(\mathbf{w})) \leq 2; \quad \mathcal{R} \geq 2 \arg \min_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}_N(\mathbf{w}) \text{ such that } k\mathbf{w}^2 k_2 \leq B_w = 2;$$

$\text{TF}^0(\mathbf{H}^{(0)})$ approximately implements (ICGD) with initialization $\mathbf{w}_{\text{GD}}^0 = \mathbf{0}$:

- (Parameter space) For every $\ell \in [L]$, the ℓ -th layer's output $\mathbf{H}^{(\ell)} = \text{TF}^0_{(1:\cdot)}(\mathbf{H}^{(0)})$ approximates ℓ steps of (ICGD): We have $\mathbf{h}_i^{(\ell)} = [\mathbf{x}_i; \mathbf{y}_i; \mathbf{w}_i^d; \mathbf{0}_{D-2d-3}; 1; t_i]$ for every $i \in [N + 1]$, where

$$\|\mathbf{w}_i^d - \mathbf{w}_{\text{GD}}^d\|_2 \leq \ell \cdot (B_x):$$

Note that the bound scales as $O(\ell)$, a linear error accumulation.

- (Prediction space) The final output $\mathbf{H}^{(L+1)} = \text{TF}^0(\mathbf{H}^{(0)})$ approximates the prediction of L steps of (ICGD): We have $\mathbf{h}_{N+1}^{(L+1)} = [\mathbf{x}_{N+1}; \mathbf{y}_{N+1}; \mathbf{w}_{N+1}^d; \mathbf{0}_{D-2d-3}; 1; t_i]$, where $\mathbf{y}_{N+1} = \mathbf{w}_{N+1}^d; \mathbf{x}_{N+1}$ so that

$$\|\mathbf{y}_{N+1} - \mathbf{w}_{\text{GD}}^d; \mathbf{x}_{N+1}\|_2 \leq (L B_x^2):$$

Further, the transformer admits norm bound $\|\mathbf{H}^{(\ell)}\|_2 \leq 2 + R + 2 C$.

The proof can be found in Appendix E.2. Theorem D.1 substantially generalizes that of von Oswald et al. [86] (which only does GD on square losses with a linear self-attention), and is simpler than the ones in Akyurek et al. [2] and Giannou et al. [32]. See Figure 4 for a pictorial illustration of the basic component of the construction, which implements a single step of gradient descent using a single attention layer (Proposition E.1).

Technically, we utilize the stability of convex gradient descent as in the following lemma (proof in Appendix E.3) to obtain the linear error accumulation in Theorem D.1; the error accumulation will become exponential in L in the non-convex case in general; see Lemma D.3(b).

Lemma D.1 (Composition of error for approximating convex GD). Suppose $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function. Let $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$, $R = 2k\mathbf{w}^* k_2$, and assume that f is L_f -smooth on $B_{2R}^d(\mathbf{w}^*)$. Let sequences $\{\mathbf{w}_g\}_{g \geq 0} \subset \mathbb{R}^d$ and $\{\mathbf{w}_{\text{GD}}^g\}_{g \geq 0} \subset \mathbb{R}^d$ be given by $\mathbf{w}_g^0 = \mathbf{w}_{\text{GD}}^0 = \mathbf{0}$,

$$\mathbf{w}_g^{+1} = \mathbf{w}_g - \frac{1}{L_f} \nabla f(\mathbf{w}_g); \quad k\mathbf{w}_g k_2 \leq R;$$

$$\mathbf{w}_{\text{GD}}^{+1} = \mathbf{w}_{\text{GD}} - \frac{1}{L_f} \nabla f(\mathbf{w}_{\text{GD}});$$

for all $g \geq 0$. Then as long as $L = L_f$, for any $0 \leq L \leq R/(2C)$, it holds that $\|\mathbf{w}_g^L - \mathbf{w}_{\text{GD}}^L\|_2 \leq L \cdot (R/2 + L) \cdot R$.

D.2 Proximal gradient descent for regularized convex losses

Proximal gradient descent (PGD) is a variant of gradient descent that is suitable for minimizing regularized risks [66], in particular those with a non-smooth regularizer such as the ℓ_1 norm. In this section, we show that transformers can approximate PGD with similar quantitative guarantees as for GD in Appendix D.1.

Let $\ell(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ be a loss function. Let $\mathcal{L}_N(\mathbf{w}) := \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{w}^\top \mathbf{x}_i; y_i) + R(\mathbf{w})$ denote the regularized empirical risk with loss function ℓ on dataset $\{(\mathbf{x}_i; y_i)\}_{i \in [N]}$ and regularizer R .

To minimize \mathcal{L}_N , we consider the proximal gradient descent trajectory on \mathcal{L}_N with initialization $\mathbf{w}_{\text{GD}}^0 = \mathbf{0} \in \mathbb{R}^d$ and learning rate $\eta > 0$:

$$\mathbf{w}_{\text{PGD}}^{t+1} := \text{prox}_{\mathcal{R}}(\mathbf{w}_{\text{PGD}}^t - \eta \nabla \mathcal{L}_N(\mathbf{w}_{\text{PGD}}^t)); \quad (\text{ICPGD})$$

where we denote $\mathcal{L}_N^0(\mathbf{w}) := \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{w}^\top \mathbf{x}_i; y_i)$.

To approximate (ICPGD) by transformers, in addition to the requirement on the loss ℓ as in Theorem D.1, we additionally require the proximal operator $\text{prox}_{\mathcal{R}}(\cdot)$ to be approximable by an MLP layer (as a vector-valued analog of Definition D.1) defined as follows.

Definition D.2 (Approximability by MLP). *An operator $P : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is $(\epsilon; R; D; C)$ -approximable by MLP, if there exists a MLP $\text{MLP}_{\text{MLP}} = (\mathbf{W}_1; \mathbf{W}_2) \in \mathbb{R}^{D \times d} \times \mathbb{R}^{d \times D}$ with hidden dimension D , $\|\mathbf{W}_1\|_{\text{op}} + \|\mathbf{W}_2\|_{\text{op}} \leq C$, such that $\sup_{\|\mathbf{w}\|_2 \leq R} \|P(\mathbf{w}) - \text{MLP}_{\text{MLP}}(\mathbf{w})\|_2 \leq \epsilon$.*

The definition above captures the proximal operator $\text{prox}_{\mathcal{R}}$ for a broad class of regularizers, such as the (commonly-used) L_1 and L_2 regularizer listed in the following proposition, for all of which one can directly check that they can be exactly implemented by an MLP as stated below.

Proposition D.1 (Proximal operators for commonly-used regularizers). *For regularizer R in $\mathcal{F}(k_1; \frac{1}{2}k_2^2; \text{Id}_{B_1(B)})$, the operator $\text{prox}_{\mathcal{R}} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is exactly approximable by MLP. More concretely, we have*

1. For $R = k_1 \|\cdot\|_1$, $\text{prox}_{\mathcal{R}}$ is $(\epsilon; \frac{1}{2}k_1; 4d; 4 + 2\epsilon)$ -approximable by MLP.
2. For $R = \frac{1}{2}k_2 \|\cdot\|_2^2$, $\text{prox}_{\mathcal{R}}$ is $(\epsilon; \frac{1}{2}k_2; 2d; 2 + 2\epsilon)$ -approximable by MLP.
3. For $R = \text{Id}_{B_1(B)}$, $\text{prox}_{\mathcal{R}} = \text{Proj}_{B_1(B)}$ is $(\epsilon; \frac{1}{2}k_2; 2d; 2 + 2\epsilon)$ -approximable by MLP.

Theorem D.2 (Convex ICPGD). *Fix any $B_w > 0$, $L > 1$, $\epsilon > 0$, and $\epsilon + \epsilon' \leq B_w/(2L)$. Suppose that*

1. The loss $\ell(\cdot)$ is convex in the first argument;
2. ℓ is $(\epsilon; R; M; C)$ -approximable by sum of relus with $R = \max\{B_x B_w; B_y\} \eta$.
3. R convex, and the proximal operator $\text{prox}_{\mathcal{R}}(\mathbf{w})$ is $(\epsilon'; R'; D'; C')$ -approximable by MLP with $R' = \sup_{\|\mathbf{w}\|_2 \leq B_w} \|\text{prox}_{\mathcal{R}}(\mathbf{w})\|_2 + \epsilon'$.

Then there exists a transformer TF with $(L + 1)$ layers, $\max_{\ell \in [L]} M^{(\ell)} = M$ heads within the first L layers, $M^{(L+1)} = 2$, and hidden dimension D' such that, for any input data $(D; \mathbf{x}_{N+1})$ such that

$$\sup_{\|\mathbf{w}\|_2 \leq B_w} \max_{\ell \in [L]} (r^{\ell} \mathcal{L}_N(\mathbf{w})) \leq \epsilon; \quad \exists \mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}_N(\mathbf{w}) \text{ such that } \|\mathbf{w}^*\|_2 \leq B_w/2;$$

TF $(\mathbf{H}^{(0)})$ approximately implements (ICGD):

1. (Parameter space) For every $\ell \in [L]$, the ℓ -th layer's output $\mathbf{H}^{(\ell)} = \text{TF}_{(\ell; \cdot)}(\mathbf{H}^{(0)})$ approximates ℓ steps of (ICGD): We have $\mathbf{h}_i^{(\ell)} = [\mathbf{x}_i; y_i; \mathbf{w}^\ell; \mathbf{0}_{D-2d-3}; 1; t_i]$ for every $i \in [N + 1]$, where

$$\mathbf{w}^\ell = \mathbf{w}_{\text{PGD}}^\ell - \eta \nabla \mathcal{L}_N(\mathbf{w}_{\text{PGD}}^\ell) \quad (\epsilon + \epsilon') (L B_x):$$

2. (Prediction space) The final output $\mathbf{H}^{(L+1)} = \text{TF}_{(L+1; \cdot)}(\mathbf{H}^{(0)})$ approximates the prediction of L steps of (ICGD): We have $\mathbf{h}_{N+1}^{(L+1)} = [\mathbf{x}_{N+1}; y_{N+1}; \mathbf{w}^L; \mathbf{0}_{D-2d-3}; 1; t_i]$, where $y_{N+1} = \mathbf{w}^L \cdot \mathbf{x}_{N+1}$ so that

$$y_{N+1} = \mathbf{w}_{\text{PGD}}^L \cdot \mathbf{x}_{N+1} \quad (\epsilon + \epsilon') (2L B_x^2):$$

Further, the weight matrices have norm bounds $\|W\| \leq 3 + R + 2C + C'$.

The proof of Theorem D.2 is essentially similar to the proof of Theorem D.1, using the following generalized version of Lemma D.1.

Lemma D.2 (Composition of error for approximating convex PGD). *Suppose $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function and R is a convex regularizer. Let $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + R(\mathbf{w})$, $R \leq 2k\mathbf{w}^*k_2$, and assume that $r \circ f$ is L_f -smooth on $B_2^d(R)$. Let sequences $\{\mathbf{w}^l\}_{l \geq 0} \subset \mathbb{R}^d$ and $\{\mathbf{w}_{\text{GD}}^l\}_{l \geq 0} \subset \mathbb{R}^d$ be given by $\mathbf{w}^0 = \mathbf{w}_{\text{GD}}^0 = \mathbf{0}$,*

$$\begin{cases} \mathbf{w}^{l+1} = \text{prox}_R(\mathbf{w}^l; r \circ f(\mathbf{w}^l) + \epsilon^l); & k\mathbf{w}^l k_2 \leq \epsilon^l; \\ \mathbf{w}_{\text{GD}}^{l+1} = \text{prox}_R(\mathbf{w}_{\text{GD}}^l; r \circ f(\mathbf{w}_{\text{GD}}^l)); \end{cases}$$

for all $l \geq 0$. Then as long as $\epsilon^l \leq L_f^{-1}$, for any $0 \leq L \leq R/(2\epsilon^l)$, it holds that $\|\mathbf{w}^L - \mathbf{w}_{\text{GD}}^L\|_2 \leq L\epsilon^l$ and $k\mathbf{w}^L k_2 \leq \frac{R}{2} + L\epsilon^l \leq R$.

The proof of the above lemma is done by utilizing the non-expansiveness of the PGD operator $\mathbf{w} \mapsto \text{prox}_R(\mathbf{w}; r \circ f(\mathbf{w}))$ and otherwise following the same arguments as for Lemma D.1.

D.3 Gradient descent on two-layer neural networks

We now move beyond the convex setting by showing that transformers can implement gradient descent on two-layer neural networks in context.

Suppose that the prediction function $\text{pred}(\mathbf{x}; \mathbf{w}) := \sum_{k=1}^K u_k r(\mathbf{v}_k^\top \mathbf{x})$ is given by a two-layer neural network, parameterized by $\mathbf{w} = [\mathbf{v}_k; u_k]_{k \in [K]} \in \mathbb{R}^{K(d+1)}$. Consider the empirical risk minimization problem:

$$\min_{\mathbf{w} \in \mathcal{W}} \mathcal{L}_N(\mathbf{w}) := \frac{1}{2N} \sum_{i=1}^N \ell(\text{pred}(\mathbf{x}_i; \mathbf{w}); y_i) = \frac{1}{2N} \sum_{i=1}^N \sum_{k=1}^K u_k r(\mathbf{v}_k^\top \mathbf{x}_i); \quad (13)$$

where \mathcal{W} is a bounded domain. For the sake of simplicity, in the following discussion we assume that $\text{Proj}_{\mathcal{W}}$ can be exactly implemented by a MLP layer (e.g. $\mathcal{W} = B_\infty(R_W)$ for some $R_W > 0$).

Theorem D.3 (Approximate ICGD on two-layer NNs). *Fix any $B_V, B_U > 0$, $L \geq 1$, $\epsilon > 0$, and $\epsilon' > 0$. Suppose that*

1. Both the activation function r and the loss function ℓ is C^4 -smooth;
2. \mathcal{W} is a closed domain such that $\mathcal{W} \cap B_\infty(R) = \{[\mathbf{v}_k; u_k]_{k \in [K]} \in \mathbb{R}^{K(d+1)} : k\mathbf{v}_k k_2 \leq B_V; |u_k| \leq B_U\}$, and $\text{Proj}_{\mathcal{W}} = \text{MLP}_{\text{mlp}}$ for some MLP layer mlp with hidden dimension D_W and $\|\text{mlp}\| \leq C_W$;

Then there exists a $(2L)$ -layer transformer TF with

$$\max_{l \in [2L]} M^{(l)} \leq \epsilon^{-2}; \quad \max_{l \in [2L]} D^{(l)} \leq \epsilon^{-2} + D_W; \quad \|W\| \leq O(1 + \epsilon) + C_W;$$

where $O(\cdot)$ hides the constants that depend on K , the radius parameters B_X, B_Y, B_U, B_V and the smoothness of r and ℓ , such that for any input data $(D; \mathbf{x}_{N+1})$ such that input sequence $\mathbf{H}^{(0)} \in \mathbb{R}^{D \times (N+1)}$ takes form (3), $\text{TF}(\mathbf{H}^{(0)})$ approximately implements in-context gradient descent on risk (13): For every $l \in [2L]$, the $2l$ -th layer's output $\mathbf{h}_i^{(2l)} = [\mathbf{x}_i; y_i; \mathbf{w}^l; \mathbf{0}; 1; t_i]$ for every $i \in [N+1]$, and

$$\mathbf{w}^l = \text{Proj}_{\mathcal{W}}(\mathbf{w}^{l-1} - (r \circ \mathcal{L}_N(\mathbf{w}^{l-1}) + \epsilon^{l-1})); \quad \mathbf{w}^0 = \mathbf{0}; \quad (14)$$

where $\epsilon^{l-1} \leq \epsilon'$ is an error term.

As a direct corollary, the transformer constructed above can approximate the true gradient descent trajectory $\{\mathbf{w}_{\text{GD}}^l\}_{l \geq 0}$ on (16), defined as $\mathbf{w}_{\text{GD}}^0 = \mathbf{0}$ and $\mathbf{w}_{\text{GD}}^{l+1} = \text{prox}_R(\mathbf{w}_{\text{GD}}^l; r \circ \mathcal{L}_N(\mathbf{w}_{\text{GD}}^l))$ for all $l \geq 0$.

Corollary D.1 (Approximating multi-step ICGD on two-layer NNs). *For any $L \geq 1$, under the same setting as Theorem D.3, the $(2L)$ -layer transformer TF there approximates the true gradient descent trajectory $\{\mathbf{w}_{\text{GD}}^g\}_{g \geq 0}$. For the intermediate iterates $\{\mathbf{w}^g\}_{g \in [L]}$ considered therein, we have*

$$\|\mathbf{w}^g - \mathbf{w}_{\text{GD}}^g\|_2 \leq L^{-1}(1 + L_f)^g;$$

where $L_f = \sup_{\mathbf{w} \in \mathcal{W}} \|\nabla^2 \mathcal{L}_N(\mathbf{w})\|_{\text{op}}$ denotes the smoothness of \mathcal{L}_N within \mathcal{W} .

Remark on error accumulation Note that in Corollary D.1, the error accumulates *exponentially* in g rather than linearly as in Theorem D.1. This is as expected, since gradient descent on non-convex objectives is inherently unstable at a high level (a slight error added upon each step may result in a drastically different trajectory); technically, this happens as the stability-like property Lemma D.1 no longer holds for the non-convex case.

Corollary D.1 is a simple implication of Theorem D.3 and Part (b) of the following convergence and trajectory closeness result for inexact gradient descent. For any closed convex set $\mathcal{W} \subseteq \mathbb{R}^d$, any function $f: \mathcal{W} \rightarrow \mathbb{R}$, and any initial point $\mathbf{w} \in \mathcal{W}$, let

$$G_{\mathcal{W}}^f(\mathbf{w}) := \frac{\mathbf{w} - \text{Proj}_{\mathcal{W}}(\mathbf{w} - \eta \nabla f(\mathbf{w}))}{\eta}$$

denote the gradient mapping at \mathbf{w} with step size η , a standard measure of stationarity in constrained optimization [63]. Note that $G_{\mathcal{W}}^f(\mathbf{w}) = \nabla f(\mathbf{w})$ when $\mathbf{w} - \eta \nabla f(\mathbf{w}) \in \mathcal{W}$ (so that the projection does not take effect).

Lemma D.3 (Convergence and trajectory closeness of inexact GD). *Suppose $f: \mathcal{W} \rightarrow \mathbb{R}$, where $\mathcal{W} \subseteq \mathbb{R}^d$ is a convex closed domain and ∇f is L_f -Lipschitz on \mathcal{W} . Let sequence $\{\mathbf{w}^g\}_{g \geq 0} \subseteq \mathbb{R}^d$ be given by $\mathbf{w}^0 = \mathbf{w}^0$,*

$$\mathbf{w}^{g+1} = \text{Proj}_{\mathcal{W}}(\mathbf{w}^g - \eta(\nabla f(\mathbf{w}^g) + \mathbf{u}^g)); \quad \|\mathbf{u}^g\|_2 \leq \epsilon$$

for all $g \geq 0$. Then the following holds.

(a) As long as $\epsilon \leq 1/L_f$, for all $L \geq 1$,

$$\min_{g \in [L-1]} \|G_{\mathcal{W}}^f(\mathbf{w}^g)\|_2^2 \leq \frac{1}{L} \sum_{g=0}^{L-1} \|G_{\mathcal{W}}^f(\mathbf{w}^g)\|_2^2 \leq \frac{8(f(\mathbf{w}^0) - \inf_{\mathbf{w} \in \mathcal{W}} f(\mathbf{w}))}{L} + 10\epsilon^2.$$

(b) Let the sequences $\{\mathbf{w}_{\text{GD}}^g\}_{g \geq 0} \subseteq \mathbb{R}^d$ and be given by $\mathbf{w}_{\text{GD}}^0 = \mathbf{w}^0$ and $\mathbf{w}_{\text{GD}}^{g+1} = \text{Proj}_{\mathcal{W}}(\mathbf{w}_{\text{GD}}^g - \eta \nabla f(\mathbf{w}_{\text{GD}}^g))$. Then it holds that

$$\|\mathbf{w}^g - \mathbf{w}_{\text{GD}}^g\|_2 \leq L_f^{-1}(1 + L_f)^g \epsilon; \quad \forall g \geq 0.$$

E Proofs for Section D

E.1 Approximating a single GD step

Proposition E.1 (Approximating a single GD step by a single attention layer). *Let $\psi(\cdot; \cdot): \mathbb{R}^2 \rightarrow \mathbb{R}$ be a loss function such that $\psi(\cdot; \cdot)$ is $(\epsilon; R; M; C)$ -approximable by sum of relus with $R = \max\{B_x, B_w; B_y; 1\}$. Let $\mathcal{L}_N(\mathbf{w}) := \frac{1}{N} \sum_{i=1}^N \psi(\mathbf{w}^\top \mathbf{x}_i; y_i)$ denote the empirical risk with loss function ψ on dataset $\{(\mathbf{x}_i; y_i)\}_{i \in [N]}$.*

Then, for any $\epsilon > 0$, there exists an attention layer $\text{Attn} = f(\mathbf{Q}_m; \mathbf{K}_m; \mathbf{V}_m)g_{m \in [M]}$ with M heads such that, for any input sequence that takes form $\mathbf{h}_i = [\mathbf{x}_i; y_i; \mathbf{w}; \mathbf{0}_{D-2d-3}; 1; t_i]$ with $\|\mathbf{w}\|_2 \leq B_w$, it gives output $\hat{\mathbf{h}}_i = [\text{Attn}(\mathbf{H})]_i = [\mathbf{x}_i; y_i; \mathbf{w}; \mathbf{0}_{D-2d-3}; 1; t_i]$ for all $i \in [N+1]$, where

$$\|\mathbf{w} - \mathbf{w}_{\text{GD}}^g\|_2 \leq \epsilon \quad (\forall B_x):$$

Further, $M \leq 2 + R + 2/C$.

Proof of Proposition E.1. As σ is $(R; M; C)$ -approximable by sum of relus, there exists a function $f: [R; R]^2 \rightarrow \mathbb{R}$ of form

$$f(s; t) = \sum_{m=1}^M c_m (a_m s + b_m t + d_m) \quad \text{with} \quad \sum_{m=1}^M (|c_m| + |a_m| + |b_m| + |d_m|) \leq 1; \quad \forall m \in [M];$$

such that $\sup_{(s;t) \in [-R; R]^2} |f(s; t) - \sigma(s; t)| \leq \epsilon$.

Next, for every $m \in [M]$, we define matrices $\mathbf{Q}_m; \mathbf{K}_m; \mathbf{V}_m \in \mathbb{R}^{D \times D}$ such that

$$\mathbf{Q}_m \mathbf{h}_i = \begin{bmatrix} a_m \mathbf{w}^\top \mathbf{x}_j \\ b_m y_j \\ d_m \\ 0 \end{bmatrix}; \quad \mathbf{K}_m \mathbf{h}_j = \begin{bmatrix} \mathbf{x}_j^\top \\ y_j \\ 1 - t_j \\ 0 \end{bmatrix}; \quad \mathbf{V}_m \mathbf{h}_j = \begin{bmatrix} \frac{(N+1)c_m}{N} \\ 0 \\ \mathbf{x}_j \\ 0_{D-2d-1} \end{bmatrix}$$

for all $i; j \in [N+1]$. As the input has structure $\mathbf{h}_i = [\mathbf{x}_i; y_i; \mathbf{w}; \mathbf{0}_{D-2d-3}; 1; t_i]$, these matrices indeed exist, and further it is straightforward to check that they have norm bounds

$$\max_{m \in [M]} \|\mathbf{Q}_m\|_{\text{op}} \leq 3; \quad \max_{m \in [M]} \|\mathbf{K}_m\|_{\text{op}} \leq 2 + R; \quad \max_{m \in [M]} \|\mathbf{V}_m\|_{\text{op}} \leq C;$$

Consequently, $\sum_{m=1}^M \|\mathbf{V}_m\|_{\text{op}} \leq 2 + R + 2C$.

Now, for every $i; j \in [N+1]$, we have

$$\begin{aligned} (\mathbf{h}_i^\top \mathbf{Q}_m \mathbf{h}_i; \mathbf{K}_m \mathbf{h}_j) &= a_m \mathbf{w}^\top \mathbf{x}_j + b_m (1 - t_j) y_j + d_m \\ &= a_m \mathbf{w}^\top \mathbf{x}_j + b_m y_j + d_m \cdot 1 \cdot t_j = 1 \cdot g; \end{aligned}$$

where the last equality follows from the bound

$$a_m \mathbf{w}^\top \mathbf{x}_j + b_m (1 - t_j) y_j + d_m \leq |a_m| B_x B_w + R \leq 2R; \quad (15)$$

so that the above relu equals 0 if $t_j = 0$. Therefore,

$$\begin{aligned} & \sum_{m=1}^M c_m (\mathbf{h}_i^\top \mathbf{Q}_m \mathbf{h}_i; \mathbf{K}_m \mathbf{h}_j) \mathbf{V}_m \mathbf{h}_j \\ &= \sum_{m=1}^M c_m (a_m \mathbf{w}^\top \mathbf{x}_j + b_m y_j + d_m) \frac{(N+1)}{N} 1 \cdot t_j = 0 \cdot g [\mathbf{0}_{d+1}; \mathbf{x}_j; \mathbf{0}_2] \\ &= f(\mathbf{w}^\top \mathbf{x}_j; y_j) \frac{(N+1)}{N} 1 \cdot t_j = 0 \cdot g [\mathbf{0}_{d+1}; \mathbf{x}_j; \mathbf{0}_{D-2d-1}]; \end{aligned}$$

Thus letting the attention layer $\text{Attn}(\mathbf{H}) = f(\mathbf{V}_m; \mathbf{Q}_m; \mathbf{K}_m)_{m \in [M]}$, we have

$$\begin{aligned} \hat{\mathbf{h}}_i &= [\text{Attn}(\mathbf{H})]_i = \mathbf{h}_i + \frac{1}{N+1} \sum_{j=1}^{N+1} \sum_{m=1}^M c_m (\mathbf{h}_i^\top \mathbf{Q}_m \mathbf{h}_i; \mathbf{K}_m \mathbf{h}_j) \mathbf{V}_m \mathbf{h}_j \\ &= \mathbf{h}_i + \frac{1}{N} \sum_{j=1}^N f(\mathbf{w}^\top \mathbf{x}_j; y_j) [\mathbf{0}_{d+1}; \mathbf{x}_j; \mathbf{0}_2] \\ &= [\mathbf{x}_i; y_i; \mathbf{w}; 1; t_i] + \frac{1}{N} \sum_{j=1}^N \underbrace{\sigma(\mathbf{w}^\top \mathbf{x}_j; y_j) [\mathbf{0}_{d+1}; \mathbf{x}_j; \mathbf{0}_{D-2d-1}]}_{\substack{\text{error vector } \mathbf{z} \\ \|\mathbf{z}\| \leq \epsilon}} + [\mathbf{0}_{d+1}; \mathbf{w}; \mathbf{0}_{D-2d-1}] \\ &= [\mathbf{x}_i; y_i; \mathbf{w} + \mathbf{z}; \mathbf{0}_{D-2d-3}; 1; t_i]; \end{aligned}$$

where the error vector $\mathbf{z} \in \mathbb{R}^d$ satisfies

$$\|\mathbf{z}\|_2 \leq \frac{1}{N} \sum_{j=1}^N |f(\mathbf{w}^\top \mathbf{x}_j; y_j) - \sigma(\mathbf{w}^\top \mathbf{x}_j; y_j)| \|\mathbf{x}_j\|_2$$

$$\frac{1}{N} \sum_{j=1}^N f(\mathbf{w}^\top \mathbf{x}_j; y_j) \approx_s \mathbb{E} f(\mathbf{w}^\top \mathbf{x}_j; y_j) \quad \|\mathbf{x}_j\|_2 \leq B_x$$

$$\frac{1}{N} \sum_{j=1}^N \|\mathbf{x}_j\|_2 \leq B_x =: B_x$$

This is the desired result. \square

E.2 Proof of Theorem D.1

We first prove part (a), which requires constructing the first L layers of \mathcal{F} . Note that by our precondition $L \leq B_w = (2^L)$.

By our precondition, the partial derivative of the loss \mathcal{L}_s is $(\cdot; R; M; C)$ -approximable by sum of relus. Therefore we can apply Proposition E.1 to obtain that, there exists a single attention layer $\text{Attn}^{(1)} = f(\mathbf{Q}_m; \mathbf{K}_m; \mathbf{V}_m)_{g_{m \in [M]}}$ with M heads (and norm bounds specified in Proposition E.1), such that for any \mathbf{w} with $\|\mathbf{w}\|_2 \leq B_w$, the attention layer $\text{Attn}^{(1)}$ maps the input $\mathbf{h}_i = [\mathbf{x}_i; \mathbf{y}_i; \mathbf{w}; \mathbf{0}_{D-2d-3}; 1; t_i]$ to output $\mathbf{h}_i' = [\mathbf{x}_i; \mathbf{y}_i; \mathbf{w}; \mathbf{0}_{D-2d-3}; 1; t_i]$ for all $i \in [N+1]$, where

$$\|\mathbf{w}\|_2 \leq B_w \implies \|\mathbf{h}_i\|_2 \leq B_x =: B_x$$

Consider the L -layer transformer $\mathcal{F}^{1:L} = (\text{Attn}^{(1)}; \dots; \text{Attn}^{(1)})$ which stacks the same attention layer $\text{Attn}^{(1)}$ for L times, and for the given input $\mathbf{h}_i^{(0)} = [\mathbf{x}_i; \mathbf{y}_i; \mathbf{w}^0; \mathbf{0}_{D-2d-3}; 1; t_i]$, its ℓ -th layer's output $\mathbf{h}_i^{(\ell)} = [\mathbf{x}_i; \mathbf{y}_i; \mathbf{w}; \mathbf{0}_{D-2d-3}; 1; t_i]$.

We now inductively show that $\|\mathbf{w}^{(\ell)}\|_2 \leq B_w$ and $\|\mathbf{w}_{\text{GD}}^{(\ell)}\|_2 \leq B_w$ for all $\ell \in [L]$. The base case of $\ell = 0$ is trivial. Suppose the claim holds for ℓ . Then for $\ell + 1 \leq L$ ($B_w = (2^L)$), the sequence $f(\mathbf{w}^{(\ell)}; g_{i \leq \ell+1})$ and $f(\mathbf{w}_{\text{GD}}^{(\ell)}; g_{i \leq \ell+1})$ satisfies the precondition of the error composition lemma (Lemma D.1) with error bound ϵ , from which we obtain $\|\mathbf{w}^{(\ell+1)}\|_2 \leq B_w$ and

$$\|\mathbf{w}_{\text{GD}}^{(\ell+1)}\|_2 \leq (B_w + \epsilon)$$

This finishes the induction, and gives the following approximation guarantee for all $\ell \in [L]$:

$$\|\mathbf{w}^{(\ell)} - \mathbf{w}_{\text{GD}}^{(\ell)}\|_2 \leq \epsilon \cdot (L B_x)$$

which proves part (a).

We now prove part (b), which requires constructing the last attention layer $\text{Attn}^{(L+1)}$. Recall $\mathbf{h}_i^{(L)} = [\mathbf{x}_i; \mathbf{y}_i; \mathbf{w}^L; \mathbf{0}_{D-2d-3}; 1; t_i]$ for all $i \in [N+1]$. We construct a 2-head attention layer $\text{Attn}^{(L+1)} = f(\mathbf{Q}_m^{(L+1)}; \mathbf{K}_m^{(L+1)}; \mathbf{V}_m^{(L+1)})_{g_{m=1,2}}$ such that for every $i, j \in [N+1]$,

$$\mathbf{Q}_1^{(L+1)} \mathbf{h}_i^{(L)} = [\mathbf{x}_i; \mathbf{0}_{D-d}]; \quad \mathbf{K}_1^{(L+1)} \mathbf{h}_j^{(L)} = [\mathbf{w}^L; \mathbf{0}_{D-d}]; \quad \mathbf{V}_1^{(L+1)} \mathbf{h}_j^{(L)} = [0_d; 1; \mathbf{0}_{D-d-1}];$$

$$\mathbf{Q}_2^{(L+1)} \mathbf{h}_i^{(L)} = [\mathbf{x}_i; \mathbf{0}_{D-d}]; \quad \mathbf{K}_2^{(L+1)} \mathbf{h}_j^{(L)} = [\mathbf{w}^L; \mathbf{0}_{D-d}]; \quad \mathbf{V}_2^{(L+1)} \mathbf{h}_j^{(L)} = [0_d; 1; \mathbf{0}_{D-d-1}];$$

Note that the weight matrices have norm bound

$$\max_{i=1,2} \|\mathbf{Q}_i^{(L+1)}\|_{\text{op}} \leq 1; \quad \max_{i=1,2} \|\mathbf{K}_i^{(L+1)}\|_{\text{op}} \leq 1; \quad \sum_{i=1}^2 \|\mathbf{V}_i^{(L+1)}\|_{\text{op}} \leq 2;$$

Then we have

$$\begin{aligned} \mathbf{h}_{N+1}^{(L+1)} &= \mathbf{h}_{N+1}^{(L)} + \frac{1}{N+1} \sum_{j=1}^{N+1} \sum_{m=1}^2 \mathbf{Q}_m^{(L+1)} \mathbf{h}_{N+1}^{(L)}; \mathbf{K}_m^{(L+1)} \mathbf{h}_j^{(L)} \mathbf{V}_m^{(L+1)} \mathbf{h}_j^{(L)} \\ &= [\mathbf{x}_i; \mathbf{0}; \mathbf{w}^L; \mathbf{0}_{D-2d-3}; 1; 1] + \left(\sum_{j=1}^N \mathbf{w}^L; \mathbf{x}_{N+1} \right) \left(\sum_{j=1}^N \mathbf{w}^L; \mathbf{x}_{N+1} \right) [0_d; 1; \mathbf{0}_{D-d-1}] \\ &\stackrel{(i)}{=} [\mathbf{x}_i; \mathbf{0}; \mathbf{w}^L; \mathbf{0}_{D-2d-3}; 1; 1] + [0_d; \mathbf{w}^L; \mathbf{x}_{N+1}; \mathbf{0}_{D-d-1}] \\ &= [\mathbf{x}_i; \underbrace{\mathbf{w}^L; \mathbf{x}_{N+1}}_{\mathbf{z}_{N+1}}; \mathbf{w}^L; \mathbf{0}_{D-2d-3}; 1; 1]; \end{aligned}$$

1. (Parameter space) For every $\ell \in [L]$, the ℓ -th layer's output $\mathbf{H}^{(\ell)} = \text{TF}_{(\ell)}(\mathbf{H}^{(0)})$ approximates ℓ steps of (ICGD-2): We have $\mathbf{h}_i^{(\ell)} = [\mathbf{x}_i; \mathbf{y}_i; \mathbf{w}; \mathbf{0}_{D-2d-3}; 1; t_i]$ for every $i \in [N+1]$, where

$$\mathbf{w} = \mathbf{w}_{\text{GD}} \quad (2L B_x):$$

2. (Prediction space) The final output $\mathbf{H}^{(L+1)} = \text{TF}(\mathbf{H}^{(0)})$ approximates the prediction of L steps of (ICGD-2): We have $\mathbf{h}_{N+1}^{(L+1)} = [\mathbf{x}_{N+1}; \mathbf{y}_{N+1}; \mathbf{w}^L; \mathbf{0}_{D-2d-3}; 1; 0]$, where

$$\mathbf{y}_{N+1} = \mathbf{w}_{\text{GD}}^L \mathbf{x}_{N+1} \quad (2L B_x^2):$$

Further, the transformer admits norm bound $\|\cdot\| \leq 2 + R + (2C + \dots)$:

Proof. This construction is the same as in the proof of Theorem D.1, except that within each layer $\ell \in [L]$, we add one more attention head $(\mathbf{Q}^{(\ell)}; \mathbf{K}^{(\ell)}; \mathbf{V}^{(\ell)}) \in \mathbb{R}^{D \times D}$ which when acting on its input $\mathbf{h}_i^{(\ell-1)} = [\cdot; \cdot; \mathbf{w}^{(\ell-1)}; 1; \cdot]$ gives

$$\mathbf{Q}^{(\ell)} \mathbf{h}_i^{(\ell-1)} = \begin{bmatrix} 1 \\ \mathbf{0}_{D-1} \end{bmatrix}; \quad \mathbf{K}^{(\ell)} \mathbf{h}_j^{(\ell-1)} = \begin{bmatrix} 1 \\ \mathbf{0}_{D-1} \end{bmatrix}; \quad \mathbf{V}^{(\ell)} \mathbf{h}_j^{(\ell-1)} = \begin{bmatrix} 2 \\ \mathbf{w}^{(\ell-1)} \\ \mathbf{0}_2 \end{bmatrix}$$

for all $i, j \in [N+1]$. Note that $\mathbf{Q}^{(\ell)}_{\text{op}} = \mathbf{K}^{(\ell)}_{\text{op}} = 1$, and $\mathbf{V}^{(\ell)}_{\text{op}} = \dots$. Further, it is straightforward to check that the output of this attention head on every $\mathbf{h}_j^{(\ell-1)}$ is

$$\frac{1}{N+1} \sum_{j=1}^{N+1} \mathbf{Q}^{(\ell)} \mathbf{h}_i^{(\ell-1)}; \mathbf{K}^{(\ell)} \mathbf{h}_j^{(\ell-1)} \mathbf{V}^{(\ell)} \mathbf{h}_j^{(\ell-1)} = \begin{bmatrix} 2 \\ \mathbf{w}^{(\ell-1)} \\ \mathbf{0}_2 \end{bmatrix}$$

Adding this onto the original output of the ℓ -th layer exactly implements the gradient of the regularizer $\mathbf{w} \nabla_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|_2^2$. The rest of the proof follows by repeating the argument of Theorem D.1, and combining the norm bound for the additional attention head here with the norm bound therein. \square

E.5 Proof of Theorem D.3

We only need to prove the following single-step version of Theorem D.3.

Proposition E.2. Under the assumptions of Theorem D.3, there exists a 2-layer transformer TF with the same bounds on the number of heads, hidden dimension and the norm, such that for any input data $(D; \mathbf{x}_{N+1})$ and any $\mathbf{w} \in \mathbb{R}^d$, TF maps

$$\mathbf{h}_i = [\mathbf{x}_i; \mathbf{y}_i; \mathbf{w}; \mathbf{0}; 1; t_i] \quad \rightarrow \quad \mathbf{h}'_i = [\mathbf{x}_i; \mathbf{y}_i; \mathbf{w}^+; \mathbf{0}; 1; t_i];$$

where

$$\mathbf{w}^+ = \text{Proj}_{\mathcal{W}} \left(\mathbf{w} + \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} \ell(\mathbf{x}_i; \mathbf{y}_i) \right); \quad \|\mathbf{w}^+\|_2 \leq k \|\mathbf{w}\|_2 + \dots;$$

Before we present the formal (and technical) proof of Proposition E.2, we first provide some intuitions. To begin with, we first note that

$$\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} \ell(\mathbf{x}_i; \mathbf{y}_i) = \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} (\text{pred}(\mathbf{x}_i; \mathbf{w}) - \mathbf{y}_i) \quad (16)$$

where $\nabla_{\mathbf{w}}$ is the partial derivative of ℓ with respect to the first component, and

$$\nabla_{\mathbf{w}} \text{pred}(\mathbf{x}_i; \mathbf{w}) = \begin{bmatrix} \frac{\partial}{\partial w_1} r(\mathbf{w}; \mathbf{x}_i) \\ \vdots \\ \frac{\partial}{\partial w_K} r(\mathbf{w}; \mathbf{x}_i) \end{bmatrix} \in \mathbb{R}^{K(d+1)}; \quad (17)$$

Therefore, the basic idea is that we can use an attention layer to approximate $(\mathbf{x}_i; \mathbf{w}) \triangleright \text{pred}(\mathbf{x}_i; \mathbf{w})$, then use an MLP layer to implement $(\text{pred}(\mathbf{x}_i; \mathbf{w}); y'_i; t_i) \triangleright 1f_i < N + 1g \text{ @}_1 \text{ `}(\text{pred}(\mathbf{x}_i; \mathbf{w}); y_i)$, and then use an attention layer to compute the gradient descent step $\mathbf{w} \triangleright \mathbf{w} \text{ } r L_N(\mathbf{w})$, and finally use an MLP layer to implement the projection into W .

Based on the observations above, we now present the proof of Proposition E.2.

Proof of Proposition E.2. We write $D_0 = d + 1 + K(d + 1)$ be the length of the vector $[\mathbf{x}_i; y_i; \mathbf{w}]$. We also define

$$B_r := \max_{|t| \leq B_x B_u} jr(t)j; \quad B_g := \max_{|t| \leq K B_r; |y| \leq B_y} j@_1 \text{ `}(t; y)j;$$

Let us fix $\epsilon_r, \epsilon_p, \epsilon \in (0, 1)$ that will be specified later in proof (see (18)). By our assumption and Proposition B.1, the following facts hold.

- (1) The function $r(t)$ is $(\epsilon_r; R_1; M_1; C_1)$ for $R_1 = \max\{B_x B_u; 1\}g, M_1 \in C_1^{\epsilon_r - 2}$, where C_1 depends only on R_1 and the C^2 -smoothness of r . Therefore, there exists

$$\tau(t) = \sum_{m=1}^M c_m^1 (\mathbf{a}_m^1; [t; 1]) \text{ with } \sum_{m=1}^M c_m^1 \leq C_1; \mathbf{a}_m^1 \in \mathbb{R}^{1; 8m} \subseteq [M_1];$$

$$\text{such that } \sup_{t \in [-R_1; R_1]} jr(t) - \tau(t)j \leq \epsilon_r.$$

- (2) The function $(t; y) \triangleright @_1 \text{ `}(t; y)$ is $(\epsilon; R_2; M_2; C_2)$ for $R_2 = \max\{K B_r; B_y; 1\}g, M_2 \in C_2^{\epsilon - 2}$, where C_2 depends only on R_2 and the C^3 -smoothness of $@_1 \text{ `}$. Therefore, there exists

$$g(t; y) = \sum_{m=1}^M c_m^2 (\mathbf{a}_m^2; [t; y; 1]) \text{ with } \sum_{m=1}^M c_m^2 \leq C_2; \mathbf{a}_m^2 \in \mathbb{R}^{1; 8m} \subseteq [M_2];$$

$$\text{such that } \sup_{(t; y) \in [-R_2; R_2]^2} jg(t; y) - @_1 \text{ `}(t; y)j \leq \epsilon.$$

- (3) The function $(s; t) \triangleright s \text{ } r'(t)$ is $(\epsilon_p; R_3; M_3; C_3)$ for $R_3 = \max\{B_x B_u; B_g B_u; 1\}g, M_3 \in C_3^{\epsilon_p - 2}$, where C_3 depends only on R_3 and the C^3 -smoothness of r' . Therefore, there exists

$$P(s; t) = \sum_{m=1}^M c_m^3 (\mathbf{a}_m^3; [s; t; 1]) \text{ with } \sum_{m=1}^M c_m^3 \leq C_3; \mathbf{a}_m^3 \in \mathbb{R}^{1; 8m} \subseteq [M_3];$$

$$\text{such that } \sup_{(s; t) \in [-R_3; R_3]^2} jP(s; t) - s \text{ } r'(t)j \leq \epsilon_p.$$

In the following, we proceed to construct the desired transformer step by step.

Step 1: construction of $\mathbf{Q}_{k; m}^{(1)}$. We consider the matrices $\mathbf{Q}_{k; m}^{(1)}, \mathbf{K}_{k; m}^{(1)}, \mathbf{V}_{k; m}^{(1)}$ $g_{k \in [K]; m \in [M_1]}$ so that for all $i; j \in [N + 1]$, we have

$$\mathbf{Q}_{k; m}^{(1)} \mathbf{h}_i = \begin{bmatrix} \mathbf{a}_m^1[1] \mathbf{x}_i \\ \mathbf{a}_m^1[2] \\ \mathbf{0} \end{bmatrix}; \quad \mathbf{K}_{k; m}^{(1)} \mathbf{h}_j = \begin{bmatrix} \mathbf{v}_k \\ 1 \\ \mathbf{0} \end{bmatrix}; \quad \mathbf{V}_{k; m}^{(1)} \mathbf{h}_j = c_m^1 u_k \mathbf{e}_{D_0 + 1};$$

As the input has structure $\mathbf{h}_i = [\mathbf{x}_i; y'_i; \mathbf{w}; \mathbf{0}; 1; t_i]$, these matrices indeed exist, and further it is straightforward to check that they have norm bounds

$$\max_{k; m} \|\mathbf{Q}_{k; m}^{(1)}\|_{\text{op}} \leq 1; \quad \max_{k; m} \|\mathbf{K}_{k; m}^{(1)}\|_{\text{op}} \leq 1; \quad \sum_{k; m} \|\mathbf{V}_{k; m}^{(1)}\|_{\text{op}} \leq C_1.$$

A simple calculation shows that

$$\sum_{m \in [M_1]; k \in [K]} \mathbf{Q}_{k; m}^{(1)} \mathbf{h}_i; \mathbf{K}_{k; m}^{(1)} \mathbf{h}_j; \mathbf{V}_{k; m}^{(1)} \mathbf{h}_j = \sum_{k=1}^K u_k \mathcal{T}(N_k; \mathbf{x}_i) \mathbf{e}_{D_0 + 1};$$

For simplicity, we denote $\overline{\text{pred}}(\mathbf{x}; \mathbf{w}) := \prod_{k=1}^K u_k \tau(\mathbf{w}_k; \mathbf{x}_i)$ in the following analysis. Thus, letting the attention layer $\text{attn}^{(1)} = f(\mathbf{V}_{k,m}^{(1)} \cdot \mathbf{Q}_{k,m}^{(1)} \cdot \mathbf{K}_{k,m}^{(1)}) g_{(k,m)}$, we have

$$\text{Attn}_{\text{attn}}^{(1)} : \mathbf{h}_i \mapsto \mathbf{h}_i^{(0.5)} = [\mathbf{x}_i; y_i; \mathbf{w}; \overline{\text{pred}}(\mathbf{x}_i; \mathbf{w}); \mathbf{0}; 1; t_j]$$

Step 2: construction of $\text{mlp}^{(1)}$. We pick matrices $\mathbf{W}_1, \mathbf{W}_2$ so that \mathbf{W}_1 maps

$$\mathbf{W}_1 \mathbf{h}_i^{(0.5)} = \begin{matrix} \mathbf{h} \\ \mathbf{a}_m^2[1] \overline{\text{pred}}(\mathbf{x}_i; \mathbf{w}) + \mathbf{a}_m^2[2] y_i + \mathbf{a}_m^2[3] R_2(1-t_j) \\ \mathbf{0} \end{matrix} \stackrel{i}{m \in [M_2]} \in \mathbb{R}^{M_2};$$

and $\mathbf{W}_2 \in \mathbb{R}^{D \times M_3}$ with entries being $(\mathbf{W}_2)_{(j,m)} = c_m^2 1 f_j = D_0 + 2g$. It is clear that $\|\mathbf{W}_1\|_{\text{op}} \leq R_2 + 1, \|\mathbf{W}_2\|_{\text{op}} \leq C_2$. Then we have

$$\begin{aligned} \mathbf{W}_2 (\mathbf{W}_1 \mathbf{h}_i^{(0.5)}) &= \prod_{m \in [M_3]} \mathbf{a}_m^2 [\overline{\text{pred}}(\mathbf{x}_i; \mathbf{w}); y_i; 1] R_2(1-t_j) c_m^2 \mathbf{e}_{D_0+2} \\ &= 1 f t_j = 1 g \cdot g(\overline{\text{pred}}(\mathbf{x}_i; \mathbf{w}); y_i) \cdot \mathbf{e}_{D_0+2} \end{aligned}$$

In the following, we abbreviate $g_i = 1 f t_j = 1 g \cdot g(\overline{\text{pred}}(\mathbf{x}_i; \mathbf{w}); y_i)$. Hence, $\text{mlp}^{(1)}$ maps

$$\text{MLP}_{\text{mlp}}^{(1)} : \mathbf{h}_i^{(0.5)} \mapsto \mathbf{h}_i^{(1)} = [\mathbf{x}_i; y_i; \mathbf{w}; \overline{\text{pred}}(\mathbf{x}_i; \mathbf{w}); g_i; \mathbf{0}; 1; t_j]$$

By the definition of the function g , for each $i \in [N]$,

$$g_i = \mathcal{G}_1(\overline{\text{pred}}(\mathbf{x}_i; \mathbf{w}); y_i) + B_u L \cdot r_i$$

where $L := \max_{|t| \leq K B_r, |y| \leq B_y} \mathcal{G}_{tt}^2(t; y)$ is the smoothness of \mathcal{G}_1 . Also, $g_{N+1} = 0$ by definition.

Step 3: construction of $\text{attn}^{(2)}$. We consider the matrices $f \mathbf{Q}_{k,1,m}^{(2)} \cdot \mathbf{K}_{k,1,m}^{(2)} \cdot \mathbf{V}_{k,1,m}^{(2)} g_{k \in [K]; m \in [M_3]}$ so that for all $i; j \in [N+1]$, we have

$$\mathbf{Q}_{k,1,m}^{(2)} \mathbf{h}_i^{(1)} = \begin{matrix} \mathbf{a}_m^3[1] u_k \\ \mathbf{a}_m^3[2] \mathbf{v}_k \\ \mathbf{a}_m^3[3] \\ \mathbf{0} \end{matrix}; \quad \mathbf{K}_{k,1,m}^{(2)} \mathbf{h}_j^{(1)} = \begin{matrix} \mathbf{g}_j \\ \mathbf{x}_j \\ \mathbf{0} \end{matrix}; \quad \mathbf{V}_{k,1,m}^{(2)} \mathbf{h}_j^{(1)} = \frac{(N+1) c_m^3}{N} \begin{matrix} \mathbf{0}_{k(d+1)} \\ \mathbf{x}_j \\ \mathbf{0} \end{matrix} \#$$

We further consider the matrices $f \mathbf{Q}_{k,2,m}^{(2)} \cdot \mathbf{K}_{k,2,m}^{(2)} \cdot \mathbf{V}_{k,2,m}^{(2)} g_{k \in [K]; m \in [M_1]}$ so that for all $i; j \in [N+1]$, we have

$$\mathbf{Q}_{k,2,m}^{(2)} \mathbf{h}_i^{(1)} = \begin{matrix} \mathbf{a}_m^1[1] \mathbf{v}_k \\ \mathbf{a}_m^1[2] \\ \mathbf{0} \end{matrix}; \quad \mathbf{K}_{k,2,m}^{(2)} \mathbf{h}_j^{(1)} = \begin{matrix} \mathbf{x}_j \\ \mathbf{1} \\ \mathbf{0} \end{matrix} \#; \quad \mathbf{V}_{k,2,m}^{(2)} \mathbf{h}_j^{(1)} = \frac{(N+1) c_m^1}{N} \begin{matrix} \mathbf{0}_{k(d+1)+d} \\ \mathbf{g}_j \\ \mathbf{0} \end{matrix} \#$$

By the structure of the input $\mathbf{h}_i^{(1)}$, these matrices indeed exist, and further it is straightforward to check that they have norm bounds

$$\max_{(k;w;m)} \|\mathbf{Q}_{k;w;m}^{(2)}\|_{\text{op}} \leq 1; \quad \max_{(k;w;m)} \|\mathbf{K}_{k;w;m}^{(2)}\|_{\text{op}} \leq 1; \quad \prod_{(k;w;m)} \|\mathbf{V}_{k;w;m}^{(2)}\|_{\text{op}} \leq 2 C_1 + 2 C_3$$

Furthermore, a simple calculation shows that

$$\mathbf{g}(\mathbf{w}) =: \frac{1}{N+1} \prod_{i=1}^{N+1} \prod_{(k;w;m)} \mathbf{Q}_{k;w;m}^{(2)} \mathbf{h}_i^{(1)} \cdot \mathbf{K}_{k;w;m}^{(2)} \mathbf{h}_j^{(1)} \cdot \mathbf{V}_{k;w;m}^{(2)} \mathbf{h}_j^{(1)} = \frac{1}{N} \prod_{j=1}^N \begin{matrix} \mathbf{0}_{d+1} \\ P(u_1 g_j; \mathbf{w}_1; \mathbf{x}_j) \mathbf{x}_j \\ \tau(\mathbf{w}_1; \mathbf{x}_j) g_j \\ \vdots \\ P(u_K g_j; \mathbf{w}_K; \mathbf{x}_j) \mathbf{x}_j \\ \tau(\mathbf{w}_K; \mathbf{x}_j) g_j \\ \mathbf{0} \end{matrix}$$

where the summation is taken over all possibilities of the tuple $(k; w; m)$, i.e. over the union of $[K] \times [M_3]$ and $[K] \times [M_1]$.

By our definition, we have $jP(s; t) = sr'(t)j \leq \rho$ for all $s; t \in [R_3; R_3]$. Therefore, for each $i \in [N], k \in [K]$,

$$jP(u_k g_j; \mathbf{h}_k; \mathbf{x}_j; i) \leq \rho_1 \cdot (\text{pred}(\mathbf{x}_j; \mathbf{w}); y_j) + u_k \cdot r'(\mathbf{h}_k; \mathbf{x}_j; i) \leq \rho + j g_j \leq \rho_1 \cdot (\text{pred}(\mathbf{x}_j; \mathbf{w}); y_j) + j u_k \leq \rho + B_u L_r (\rho + B_u L \cdot \rho);$$

where $L_r := \max_{|t| \leq B_x B_u} j r'(t) j$ is the upper bound of r' . Similarly, for each $i \in [N], k \in [K]$, we have

$$j r(\mathbf{h}_k; \mathbf{x}_j; i) \leq g_j + r(\mathbf{h}_k; \mathbf{x}_j; i) \leq \rho_1 \cdot (\text{pred}(\mathbf{x}_j; \mathbf{w}); y_j) + 2 B_g \rho_r + 2 B_r (\rho + B_u L^2 \rho_r);$$

As for the case $i = N + 1$, we have $g_{N+1} = 0$ and $jP(u_k g_{N+1}; \mathbf{h}_k; \mathbf{x}_{N+1}; i) \leq \rho$ for each $k \in [K]$ by definition. Combining these estimations and using (16) and (17), we can conclude that

$$\| \mathbf{g}(\mathbf{w}) + r \mathbf{L}_N(\mathbf{w}) \|_2 \leq \frac{\rho}{K} B_x [\rho + B_u L_r (\rho + B_u L \cdot \rho)] + 2 \frac{\rho}{K} [B_g \rho_r + B_r (\rho + B_u L \cdot \rho)];$$

Thus, to ensure $\| \mathbf{g}(\mathbf{w}) + r \mathbf{L}_N(\mathbf{w}) \|_2 \leq \epsilon$, we only need to choose $\rho; \rho_r; \rho_r$ as

$$\rho = \frac{\epsilon}{3} \frac{1}{K B_x}; \quad \rho_r = \frac{\epsilon}{9} \frac{1}{K \max\{B_r; L_r B_x B_u\}}; \quad \rho_r = \frac{\epsilon}{15} \frac{1}{K \max\{B_g; L \cdot B_r B_u; L_r L \cdot B_x B_r B_u^2\}}; \quad (18)$$

Thus, letting the attention layer $\text{attn}^{(2)} = f(\mathbf{V}_{k;w;m}^{(2)}; \mathbf{Q}_{k;w;m}^{(2)}; \mathbf{K}_{k;w;m}^{(2)}) g_{(k;w;m)}$, we have

$$\text{Attn}_{\text{attn}}^{(2)} : \mathbf{h}_i^{(1)} \mapsto \mathbf{h}_i^{(1.5)} = [\mathbf{x}_i; y_i; \mathbf{w} + \mathbf{g}(\mathbf{w}); \overline{\text{pred}(\mathbf{x}_i; \mathbf{w})}; g_i; \mathbf{0}; 1; t_i];$$

Step 4: construction of $\text{mlp}^{(2)}$. We only need to pick $\text{mlp}^{(2)}$ so that it maps

$$\mathbf{h}_i^{(1.5)} = [\mathbf{x}_i; y_i; \mathbf{w} + \mathbf{g}(\mathbf{w}); \overline{\text{pred}(\mathbf{x}_i; \mathbf{w})}; g_i; \mathbf{0}; 1; t_i] \xrightarrow{\text{MLP}_{\text{mlp}}^{(2)}} \mathbf{h}_i^{(2)} = [\mathbf{x}_i; y_i; \text{Proj}_{\mathcal{W}}(\mathbf{w} + \mathbf{g}(\mathbf{w})); \mathbf{0}; \mathbf{0}; \mathbf{0}; 1; t_i];$$

By our assumption on the map $\text{Proj}_{\mathcal{W}}$, this is easy.

Combining the four steps above and taking $\epsilon = (\frac{1}{\text{attn}}; \frac{1}{\text{mlp}}; \frac{2}{\text{attn}}; \frac{2}{\text{mlp}})$ completes the proof. \square

E.6 Proof of Lemma D.3

For every $\epsilon > 0$, define the intermediate iterates (before projection)

$$\mathbf{w}^{\epsilon+1/2} := \mathbf{w}^\epsilon + r f(\mathbf{w}^\epsilon) + \epsilon; \quad \mathbf{w}_{\text{GD}}^{\epsilon+1/2} := \mathbf{w}_{\text{GD}}^\epsilon + r f(\mathbf{w}_{\text{GD}}^\epsilon);$$

so that $\mathbf{w}^{\epsilon+1} = \text{Proj}_{\mathcal{W}}(\mathbf{w}^{\epsilon+1/2})$ and $\mathbf{w}_{\text{GD}}^{\epsilon+1} = \text{Proj}_{\mathcal{W}}(\mathbf{w}_{\text{GD}}^{\epsilon+1/2})$.

We first prove part (a). We begin by deriving a relation between $\mathbf{w}^{\epsilon+1} - \mathbf{w}^{\epsilon+1/2}$ and $G_{\mathcal{W}}^f(\mathbf{w}^{\epsilon+1/2})$.

Let $\mathbf{w}^{\epsilon+1/2} := \mathbf{w}^\epsilon + r f(\mathbf{w}^\epsilon)$ and $\mathbf{w}^{\epsilon+1} := \text{Proj}_{\mathcal{W}}(\mathbf{w}^{\epsilon+1/2})$ denote the *exact* projected gradient iterate starting from \mathbf{w}^ϵ . We have

$$\begin{aligned} \mathbf{w}^{\epsilon+1} - \mathbf{w}^{\epsilon+1/2} &\stackrel{(i)}{\leq} \frac{1}{2} \|\mathbf{w}^{\epsilon+1} - \mathbf{w}^{\epsilon+1/2}\|_2^2 - \frac{1}{2} \|\mathbf{w}^{\epsilon+1} - \mathbf{w}^{\epsilon+1/2}\|_2^2 \stackrel{(ii)}{\leq} \frac{1}{2} \|\mathbf{w}^{\epsilon+1} - \mathbf{w}^{\epsilon+1/2}\|_2^2 - \frac{1}{2} \|\mathbf{w}^{\epsilon+1/2} - \mathbf{w}^{\epsilon+1/2}\|_2^2 \\ &\stackrel{(iii)}{=} \frac{2}{2} G_{\mathcal{W}}^f(\mathbf{w}^{\epsilon+1/2})^2 - k \epsilon^2 \frac{2}{2} G_{\mathcal{W}}^f(\mathbf{w}^{\epsilon+1/2})^2 - 2 \epsilon^2. \end{aligned} \quad (19)$$

Above, (i) uses the inequality $\|ka - bk\|_2^2 \leq \frac{1}{2} k a k_2^2 + k b k_2^2$; (ii) uses the fact that projection to a convex set is a non-expansion; (iii) uses the definition of the gradient mapping.

By the L_f -smoothness of f within \mathcal{W} , we have

$$f(\mathbf{w}^{\epsilon+1}) - f(\mathbf{w}^\epsilon) \leq r f(\mathbf{w}^\epsilon); \mathbf{w}^{\epsilon+1} - \mathbf{w}^\epsilon + \frac{L_f}{2} \|\mathbf{w}^{\epsilon+1} - \mathbf{w}^\epsilon\|_2^2$$

$$\begin{aligned}
&= \frac{\|\mathbf{w}^{\ell} - \mathbf{w}^{\ell+\frac{1}{2}}\|^2}{2} + \|\mathbf{w}^{\ell+1} - \mathbf{w}^{\ell}\|^2 + \frac{L_f}{2} \|\mathbf{w}^{\ell+1} - \mathbf{w}^{\ell}\|^2 \\
&\stackrel{(i)}{=} \frac{\|\mathbf{w}^{\ell} - \mathbf{w}^{\ell+1}\|^2}{2} + \|\mathbf{w}^{\ell+1} - \mathbf{w}^{\ell}\|^2 + \|\mathbf{w}^{\ell} - \mathbf{w}^{\ell+1}\|^2 + \frac{L_f}{2} \|\mathbf{w}^{\ell+1} - \mathbf{w}^{\ell}\|^2 \\
&= \frac{1}{4} + \frac{L_f}{2} \|\mathbf{w}^{\ell+1} - \mathbf{w}^{\ell}\|^2 + \frac{1}{4} \|\mathbf{w}^{\ell+1} - \mathbf{w}^{\ell}\|^2 + \|\mathbf{w}^{\ell} - \mathbf{w}^{\ell+1}\|^2 \\
&\stackrel{(ii)}{=} \frac{1}{4} \|\mathbf{w}^{\ell+1} - \mathbf{w}^{\ell}\|^2 + \|\mathbf{w}^{\ell} - \mathbf{w}^{\ell+1}\|^2 \\
&\stackrel{(iii)}{=} \frac{1}{4} \frac{2}{2} G_{\mathcal{W};}^f(\mathbf{w}^{\ell})^2 + \|\mathbf{w}^{\ell} - \mathbf{w}^{\ell+1}\|^2 + \|\mathbf{w}^{\ell} - \mathbf{w}^{\ell+1}\|^2 \\
&\quad \frac{1}{8} G_{\mathcal{W};}^f(\mathbf{w}^{\ell})^2 + \frac{5}{4} \|\mathbf{w}^{\ell} - \mathbf{w}^{\ell+1}\|^2.
\end{aligned}$$

Above, (i) uses the property $\|\mathbf{w}^{\ell+1} - \mathbf{w}^{\ell+\frac{1}{2}}\| + \|\mathbf{w}^{\ell+1} - \mathbf{w}^{\ell}\| = 0$ of the projection $\mathbf{w}^{\ell+1} = \text{Proj}_{\mathcal{W}}(\mathbf{w}^{\ell+\frac{1}{2}})$ (using $\mathbf{w}^{\ell} \in \mathcal{W}$); (ii) uses $L_f=2$ ($\|\cdot\| = \|\cdot\|$) by our choice of $\|\cdot\| = L_f \|\cdot\|$; (iii) uses (19).

Rearranging and summing the above over $\ell = 0, \dots, L-1$, we obtain

$$\frac{1}{8} \sum_{\ell=0}^{L-1} G_{\mathcal{W};}^f(\mathbf{w}^{\ell})^2 + f(\mathbf{w}^0) - f(\mathbf{w}^L) + \frac{5}{4} L \|\mathbf{w}^0 - \mathbf{w}^L\|^2.$$

Dividing both sides by $L=8$ yields part (a).

Next, we prove part (b). Let $C := 1 + L_f$. We prove by induction that

$$\|\mathbf{w}^{\ell} - \mathbf{w}_{\text{GD}}^{\ell}\| \leq \frac{C^{\ell} - 1}{C - 1} \|\mathbf{w}^0 - \mathbf{w}_{\text{GD}}^0\| \quad (20)$$

for all $\ell \geq 0$. The base case of $\ell = 0$ follows by definition that $\mathbf{w}^0 = \mathbf{w}_{\text{GD}}^0 = \mathbf{w}^0$. Suppose the result holds for ℓ . Then for $\ell + 1$, we have

$$\begin{aligned}
&\|\mathbf{w}^{\ell+1} - \mathbf{w}_{\text{GD}}^{\ell+1}\| \stackrel{(i)}{\leq} \|\mathbf{w}^{\ell+\frac{1}{2}} - \mathbf{w}_{\text{GD}}^{\ell+\frac{1}{2}}\| = \|\mathbf{w}^{\ell} - \mathbf{w}_{\text{GD}}^{\ell}\| + \|\mathbf{w}^{\ell} - \mathbf{w}^{\ell+\frac{1}{2}}\| + \|\mathbf{w}_{\text{GD}}^{\ell} - \mathbf{w}_{\text{GD}}^{\ell+\frac{1}{2}}\| \\
&\stackrel{(ii)}{\leq} C \|\mathbf{w}^{\ell} - \mathbf{w}_{\text{GD}}^{\ell}\| + \|\mathbf{w}^{\ell} - \mathbf{w}^{\ell+\frac{1}{2}}\| \stackrel{(iii)}{\leq} C \frac{C^{\ell} - 1}{C - 1} \|\mathbf{w}^0 - \mathbf{w}_{\text{GD}}^0\| + \|\mathbf{w}^{\ell} - \mathbf{w}^{\ell+\frac{1}{2}}\| = \frac{C^{\ell+1} - 1}{C - 1} \|\mathbf{w}^0 - \mathbf{w}_{\text{GD}}^0\|.
\end{aligned}$$

Above, (i) uses again the non-expansiveness of the convex projection $\text{Proj}_{\mathcal{W}}$; (ii) uses the fact that the operator $\mathbf{w} \mapsto \mathbf{w} - \mathbf{w} - \mathbf{w}$ is $(1 + L_f) = C$ -Lipschitz; and (iii) uses the inductive hypothesis. This proves the case for $\ell + 1$ and thus finishes the induction. We can further relax (20) into

$$\|\mathbf{w}^{\ell} - \mathbf{w}_{\text{GD}}^{\ell}\| \leq \frac{C^{\ell}}{1 + L_f} \|\mathbf{w}^0 - \mathbf{w}_{\text{GD}}^0\| = L_f^{-1} (1 + L_f)^{\ell} \|\mathbf{w}^0 - \mathbf{w}_{\text{GD}}^0\|.$$

This proves part (b). \square

F Proofs for Section 3.1

F.1 Proof of Theorem 4

Fix $\epsilon > 0, 0 < \delta < 1$ with $\delta := \frac{\epsilon}{1 + \epsilon}$, and $B_W > 0$, and consider any in-context data D such that the precondition of Theorem 4 holds. Let

$$L_{\text{ridge}}(\mathbf{w}) := \frac{1}{2N} \sum_{i=1}^N (\langle \mathbf{w}; \mathbf{x}_i \rangle - y_i)^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

denote the ridge regression loss in (ICRidge), so that $\mathbf{w}_{\text{ridge}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} L_{\text{ridge}}(\mathbf{w})$. It is a standard result that $L_{\text{ridge}}(\mathbf{w}) = \frac{1}{2} \mathbf{X}^\top \mathbf{X} \mathbf{w} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$, so that L_{ridge} is $(\frac{\lambda}{2})$ -strongly convex and $(\frac{1}{2})$ -smooth over \mathbb{R}^d .

Consider the gradient descent algorithm on the ridge loss

$$\mathbf{w}_{\text{GD}}^{t+1} = \mathbf{w}_{\text{GD}}^t - \eta \nabla L_{\text{ridge}}(\mathbf{w}_{\text{GD}}^t)$$

with initialization, learning rate, and number of steps

$$\mathbf{w}_{\text{GD}}^0 := \mathbf{0}_d; \quad \eta := \frac{1}{\lambda}; \quad T := \left\lceil 2 \log \frac{B_x B_w}{2\epsilon} \right\rceil$$

By standard convergence results for strongly convex and smooth functions (Proposition B.2), we have for all $t \geq 1$ that

$$\|\mathbf{w}_{\text{GD}}^t - \mathbf{w}_{\text{ridge}}\|_2^2 \leq \exp\left(-\frac{t}{2}\right) \|\mathbf{w}_{\text{GD}}^0 - \mathbf{w}_{\text{ridge}}\|_2^2 = \exp\left(-\frac{t}{2}\right) \|\mathbf{w}_{\text{ridge}}\|_2^2$$

Further, we have

$$\|\mathbf{w}_{\text{GD}}^T - \mathbf{w}_{\text{ridge}}\|_2 \leq \exp\left(-\frac{T}{2}\right) \|\mathbf{w}_{\text{ridge}}\|_2 \leq \frac{2\epsilon}{B_x B_w} \frac{B_w}{2} \frac{\epsilon}{B_x} \quad (21)$$

It remains to construct a transformer to approximate \mathbf{w}_{GD}^T . Notice that the problem (ICRidge) corresponds to an λ -regularized ERM with the square loss $\ell(s; t) := \frac{1}{2}(s - t)^2$, whose partial derivative $\partial_s \ell(s; t) = s - t$ is exactly a sum of two relus:

$$\partial_s \ell(s; t) = \frac{1}{2}((s - t) - 2) + \frac{1}{2}((s - t) + 2)$$

In particular, this shows that $\partial_s \ell(s; t)$ is $(0; R; 2; 4)$ -approximable for any $R > 0$, in particular for $R = \max\{B_x B_w; B_y\} / \lambda$.

Therefore, we can apply Corollary E.1 with the square loss ℓ , learning rate η , regularization strength λ and accuracy parameter $\epsilon = 0$ to obtain that there exists an attention-only transformer TF^0 with $(T + 1) \leq L$ layers such that the final output $\mathbf{h}_{N+1}^{(L)} = [\mathbf{x}_{N+1}; \mathbf{y}_{N+1}]$ with

$$\mathbf{y}_{N+1} - \mathbf{w}_{\text{GD}}^T \mathbf{x}_{N+1} = \mathbf{0}; \quad (22)$$

and number of heads $M^{(i)} = 3$ for all $i \geq [L - 1]$ (can be taken as 2 in the unregularized case $\lambda = 0$ directly by Theorem D.1), and $M^{(L)} = 2$. Further, TF^0 admits norm bound $\|\mathbf{h}_{N+1}^{(L)}\|_2 \leq 2 + R + \frac{8\epsilon}{\lambda} \leq 3R + 8(\frac{1}{\lambda})^{-1} + 1 \leq 4R + 8(\frac{1}{\lambda})^{-1}$.

Combining (21) and (22), we obtain that

$$\|\mathbf{y}_{N+1} - \mathbf{w}_{\text{ridge}} \mathbf{x}_{N+1}\|_2 = \|\mathbf{w}_{\text{GD}}^T \mathbf{x}_{N+1} - \mathbf{w}_{\text{ridge}} \mathbf{x}_{N+1}\|_2 \leq \epsilon \frac{B_x}{\lambda}$$

Further, we have $\text{read}_{\mathbf{w}}(\mathbf{h}_i^T) = \mathbf{w}_{\text{GD}}^T$ for all $i \geq [N + 1]$, where $\text{read}_{\mathbf{w}}(\mathbf{h}) := \mathbf{h}_{(d+2):(2d+1)}$ (cf. Corollary E.1), so that $\|\text{read}_{\mathbf{w}}(\mathbf{h}_i^T) - \mathbf{w}_{\text{ridge}}\|_2 \leq \epsilon \frac{B_x}{\lambda}$ as shown above. This finishes the proof. \square

F.2 Statistical analysis of in-context least squares

Consider the standard least-squares algorithm A_{LS} and least-squares estimator $\mathbf{w}_{\text{LS}} \in \mathbb{R}^d$ defined as

$$A_{\text{LS}}(D)(\mathbf{x}_{N+1}) := \mathbf{h}_{\text{LS}}(\mathbf{x}_{N+1}); \quad \mathbf{w}_{\text{LS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \in \mathbb{R}^d; \quad (\text{ICLS})$$

For any distribution \mathbb{P} over $(\mathbf{x}; y) \in \mathbb{R}^d \times \mathbb{R}$ and any estimator $\mathbf{w} \in \mathbb{R}^d$, let

$$L_{\mathbb{P}}(\mathbf{w}) := \mathbb{E}_{(\mathbf{x}; y) \sim \mathbb{P}} \frac{1}{2} (\mathbf{w} \mathbf{x} - y)^2$$

denote the expected risk of \mathbf{w} over a new test example $(\mathbf{x}'; y) \sim \mathbb{P}$.

Assumption A (Well-posedness for learning linear predictors). *We say a distribution \mathbb{P} on $\mathbb{R}^d \times \mathbb{R}$ is well-posed for learning linear predictors, if $(\mathbf{x}; y) \sim \mathbb{P}$ satisfies*

- (1) $k\mathbf{x}k_2 \leq B_x$ and $\|y\| \leq B_y$ almost surely;
- (2) The covariance $\Sigma_P := \mathbb{E}_P[\mathbf{xx}^\top]$ satisfies $\min \lambda_d \leq \Sigma_P \leq \max \lambda_d$, with $0 < \min \lambda_d \leq \max \lambda_d$, and $\lambda_{\min} := \min \lambda_d = \min \lambda_d$.
- (3) The whitened vector $\Sigma_P^{-1/2} \mathbf{x}$ is K^2 -sub-Gaussian for some $K \geq 1$.
- (4) The best linear predictor $\mathbf{w}_P^? := \mathbb{E}_P[\mathbf{xx}^\top]^{-1} \mathbb{E}_P[\mathbf{x}y]$ satisfies $k\mathbf{w}_P^?k_2 \leq B_w^?$.
- (5) We have $\mathbb{E}[(y - \langle \mathbf{h}\mathbf{x}; \mathbf{w}_P^? \rangle)^2 | \mathbf{x}] \leq \sigma^2$ with probability one (over \mathbf{x}).

Further, we say P is well-posed with canonical parameters if

$$B_x = O(d); \quad B_y = O(1); \quad B_w^? = O(1); \quad O(1); \quad \max \lambda_d = O(1); \quad K = O(1); \quad (23)$$

where $O(\cdot)$ and $O(\cdot)$ only hides absolute constants.

The following result bounds the excess risk of least squares under Assumption A with a clipping operation on the predictor; the clipping allows the result to only depend on the second moment of the noise (cf. Assumption A(5)) instead of e.g. its sub-Gaussianity, and also makes the result convenient to be directly translated to a result for transformers.

Proposition F.1 (Guarantees for in-context least squares). *Suppose distribution P satisfies Assumption A. Then as long as $N \geq O(dK^4 \log(1/\delta))$, we have the following:*

- (a) The (clipped) least squares predictor achieves small expected excess risk (fast rate) over the best linear predictor: For any clipping radius $R \leq B_y$,

$$\mathbb{E}_{\mathcal{D}; \mathbf{x}_{N+1}, y_{N+1} \sim P} \frac{1}{2} (\text{clip}_R(\langle \mathbf{w}_{\text{LS}}; \mathbf{x}_{N+1} \rangle) - y_{N+1})^2 \leq \inf_{\mathbf{w} \in \mathbb{R}^d} L_P(\mathbf{w}) + O\left(\frac{R^2}{N}\right) + \frac{d^2}{N} : \quad (24)$$

- (b) We have $\mathbb{P}(E_{\text{cov}} \leq \epsilon \wedge E_w \leq \epsilon)$ $\geq 1 - \delta$, where

$$E_{\text{cov}} = E_{\text{cov}}(D) := \frac{1}{2} \lambda_d \leq \frac{1}{2} \lambda_d \leq \frac{1}{2} \lambda_d \leq \frac{1}{2} \lambda_d; \quad (25)$$

$$E_w = E_w(D) := k\mathbf{w}_{\text{LS}}k_2 \leq B_w^? + \frac{80d^2}{N \min \lambda_d}; \quad (26)$$

Proof. We first show $\mathbb{P}(E_{\text{cov}} \leq \epsilon) \geq 1 - \delta$. Let $\mathbf{b} := \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top$, and let the whitened covariance and noise variables be denoted as

$$\mathbf{x}_i = \Sigma_P^{-1/2} \mathbf{x}_i; \quad \mathbf{e}_i := \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top = \Sigma_P^{-1/2} \mathbf{b} \Sigma_P^{-1/2}.$$

Also let $z_i := y_i - \langle \mathbf{h}\mathbf{x}_i; \mathbf{w}_P^? \rangle$ denote the ‘‘noise’’ variables. Note that

$$E_{\text{cov}} = \frac{1}{2} \lambda_d \leq \frac{1}{2} \lambda_d \leq \frac{1}{2} \lambda_d \leq \frac{1}{2} \lambda_d$$

is exactly a covariance concentration of the whitened vectors $\{\mathbf{x}_i\}_{i \in [N]}$. Recall that $\mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top] = \mathbf{I}_d$, and \mathbf{x}_i are K^2 -sub-Gaussian by assumption. Therefore, we can apply [85, Theorem 4.6.1], we have with probability at least $1 - \delta$ that

$$\mathbf{e} \preceq \mathbf{I}_d + O\left(K^2 \max\left(\frac{d + \log(1/\delta)}{N}, \frac{d + \log(1/\delta)}{N}\right)\right) \mathbf{I}_d.$$

Setting $N \geq O(K^4(d + \log(1/\delta)))$ ensures that the right-hand side above is at most $\frac{3}{2} \mathbf{I}_d$, on which event we have

$$\frac{1}{2} \lambda_d \leq \frac{3}{2} \lambda_d \leq \frac{3}{2} \lambda_d \leq \frac{3}{2} \lambda_d; \quad (27)$$

i.e. E_{cov} holds. This shows that $P(E_{\text{cov}}^c) = 10^{-10}$.

Next, we show (24). Using E_{cov} , we decompose the risk as

$$\begin{aligned}
& E \frac{1}{2} (\text{clip}_R(h_{\text{wLS}}; \mathbf{x}_{N+1}) - y_{N+1})^2 \\
&= E \frac{1}{2} (\text{clip}_R(h_{\text{wLS}}; \mathbf{x}_{N+1}) - y_{N+1})^2 1_{E_{\text{cov}}} + E \frac{1}{2} (\text{clip}_R(h_{\text{wLS}}; \mathbf{x}_{N+1}) - y_{N+1})^2 1_{E_{\text{cov}}^c} \\
&\stackrel{(i)}{=} E \frac{1}{2} (h_{\text{wLS}}; \mathbf{x}_{N+1} - y_{N+1})^2 1_{E_{\text{cov}}} + 2R^2 \quad (=20) \\
&\stackrel{(ii)}{=} E_{\mathcal{D}; \mathbf{x}_{N+1}} \frac{1}{2} (h_{\text{wLS}}; \mathbf{w}_P^2; \mathbf{x}_{N+1})^2 1_{E_{\text{cov}}} + E_{\mathbf{x}_{N+1}; y_{N+1}} \frac{1}{2} (h_{\text{w}_P^2}; \mathbf{x}_{N+1} - y_{N+1})^2 1_{E_{\text{cov}}} + O(R^2) \\
&= E_{\mathcal{D}} \frac{1}{2} k_{\text{wLS}} \mathbf{w}_P^2 k_P^2 1_{E_{\text{cov}}} + E_{\mathbf{x}_{N+1}; y_{N+1}} \underbrace{\frac{1}{2} (h_{\text{w}_P^2}; \mathbf{x}_{N+1} - y_{N+1})^2}_{L_P(\mathbf{w}_P^2)} + O(R^2); \tag{28}
\end{aligned}$$

Above, (i) follows by assumption that $|y_{N+1}| \leq B_y + R$ almost surely, so that removing the clipping can only potentially increase the distance in the first term, and the square loss is upper bounded by $\frac{1}{2} (2R)^2$ almost surely in the second term; (ii) follows by the fact that $E_{\mathbf{x}_{N+1}; y_{N+1}} [h_{\text{wLS}}; \mathbf{w}_P^2; \mathbf{x}_{N+1} - (h_{\text{w}_P^2}; \mathbf{x}_{N+1} - y_{N+1})] = 0$ by the definition of \mathbf{w}_P^2 , as well as the fact that $1_{E_{\text{cov}}}$ is independent of $(\mathbf{x}_{N+1}; y_{N+1})$.

It thus remains to bound $E_{\mathcal{D}} \frac{1}{2} k_{\text{wLS}} \mathbf{w}_P^2 k_P^2 1_{E_{\text{cov}}}$. Note that on the event E_{cov} , we have

$$\frac{1}{P} \mathbf{b}^{-1} \mathbf{b}^{-1} \frac{1}{P} = \frac{1}{P} \mathbf{b}^{-1} \mathbf{b}^{-1} \frac{1}{P} = 2 \mathbf{I}_d;$$

Therefore,

$$\begin{aligned}
& \frac{1}{2} k_{\text{wLS}} \mathbf{w}_P^2 k_P^2 1_{E_{\text{cov}}} = \frac{1}{2} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \mathbf{w}_P^2 \mathbf{w}_P^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \mathbf{w}_P^2 1_{E_{\text{cov}}} \\
&= \frac{1}{2} \mathbf{z}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{P} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{z} 1_{E_{\text{cov}}} \\
&= \frac{1}{2N^2} \mathbf{z}^T \mathbf{X} \frac{1}{P} \mathbf{b}^{-1} \mathbf{b}^{-1} \frac{1}{P} \mathbf{z} 1_{E_{\text{cov}}} = \frac{1}{2N^2} \mathbf{z}^T \mathbf{X} \frac{1}{P} \mathbf{b}^{-1} \mathbf{b}^{-1} \frac{1}{P} \mathbf{z} 1_{E_{\text{cov}}} \\
&= \frac{2}{N^2} \frac{1}{P} \mathbf{z}^T \mathbf{X} \frac{1}{2} 1_{E_{\text{cov}}} = \frac{2}{N^2} \sum_{i=1}^N \mathbf{x}_i z_i 1_{E_{\text{cov}}} = \frac{2}{N^2} \sum_{i=1}^N \mathbf{x}_i z_i;
\end{aligned}$$

Note that $E[\mathbf{x}_i z_i] = \frac{1}{P} E[\mathbf{x}_i (y_i - h_{\text{w}_P^2}; \mathbf{x}_i)] = 0$. Therefore, taking expectation on the above (over \mathcal{D}), we get

$$E_{\mathcal{D}} \frac{1}{2} k_{\text{wLS}} \mathbf{w}_P^2 k_P^2 1_{E_{\text{cov}}} = \frac{2}{N^2} E \sum_{i=1}^N \mathbf{x}_i z_i^2 = \frac{2}{N} E \sum_{i=1}^N \mathbf{x}_i^T \mathbf{x}_i = \frac{2}{N} E \sum_{i=1}^N \|\mathbf{x}_i\|^2; \tag{29}$$

$$\stackrel{(i)}{=} \frac{2}{N} E \sum_{i=1}^N \|\mathbf{x}_i\|^2 = \frac{2d}{N}. \tag{30}$$

Above, (i) follows by conditioning on \mathbf{x}_1 and using Assumption A(5). Combining with (28), we obtain

$$E \frac{1}{2} (\text{clip}_R(h_{\text{wLS}}; \mathbf{x}_{N+1}) - y_{N+1})^2 \leq L_P(\mathbf{w}_P^2) + O(R^2) + \frac{d}{N};$$

This proves (24).

Finally, we show $P(E_{\text{cov}} \setminus E_w) = 10^{-10}$. Using (29) and $\frac{1}{P} \min \|\mathbf{x}_i\|^2$ by assumption, we get

$$E \frac{1}{2} k_{\text{wLS}} \mathbf{w}_P^2 k_P^2 1_{E_{\text{cov}}} \leq \frac{4d}{N \min \|\mathbf{x}_i\|^2};$$

Therefore, using an argument similar to Chebyshev's inequality,

$$\begin{aligned}
 P(E_{\text{cov}} \setminus E_w^c) &= E \left[1 - \frac{\|w_{\text{LS}}\|_2^2}{\|w\|_2^2} \right] > \frac{20}{N_{\min}} \frac{4d^2}{N_{\min}} + B_w^2 g^5 \\
 &= E \left[1 - \frac{\|w_{\text{LS}}\|_2^2}{\|w\|_2^2} \right] > \frac{20}{N_{\min}} \frac{4d^2}{N_{\min}} g^5 \\
 &= E \left[1 - \frac{\|w_{\text{LS}}\|_2^2}{\|w\|_2^2} \right] \geq \frac{20}{N_{\min}} \frac{4d^2}{N_{\min}} =: \epsilon
 \end{aligned}$$

This implies that

$$P(E_{\text{cov}} \setminus E_w) = P(E_{\text{cov}}) - P(E_{\text{cov}} \setminus E_w^c) \geq 1 - \epsilon = 1 - \epsilon =: \epsilon.$$

This is the desired result. \square

F.3 Proof of Corollary 5

The proof follows by first checking the well-conditionedness of the data D (cf. (5)) with high probability, then invoking Theorem 4 (for approximation least squares) and Proposition F.1 (for the statistical power of least squares).

First, as P satisfies Assumption A, by Proposition F.1, as long as $N \geq O(K^4(d + \log(1/\epsilon)))$, we have with probability at least $1 - \epsilon$ that event $E_{\text{cov}} \setminus E_w$ holds. On this event, we have

$$\begin{aligned}
 \frac{1}{2} \min_{d'} \lambda_{d'} &\geq \frac{1}{2} \lambda_{\min}(\mathbf{X}^T \mathbf{X}) \geq \frac{1}{2} \lambda_{\min}(\mathbf{X}^T \mathbf{X}) \geq \frac{1}{2} \lambda_{\min}(\mathbf{X}^T \mathbf{X}) \\
 \|w_{\text{LS}}\|_2 &\leq B_w := O\left(\frac{d^2}{N_{\min}}\right) + \frac{d^2}{N_{\min}} A;
 \end{aligned}$$

and thus the dataset D is well-conditioned (in the sense of (5)) with parameters $\lambda_{\min} = \frac{1}{2} \lambda_{\min}$, $\lambda_{\max} = 2 \lambda_{\max}$, and B_w defined as above. Note that the condition number of \mathbf{b} is upper bounded by $\kappa = \frac{\lambda_{\max}}{\lambda_{\min}} \leq \frac{4 \lambda_{\max}}{\lambda_{\min}}$, where κ is the upper bound on the condition number of \mathbf{p} as in Assumption A(c).

Define parameters

$$\gamma = \frac{d^2}{N}; \quad \beta = \frac{d^2}{B_y^2 N} \wedge 1; \quad (31)$$

Note that $B_w = O\left(\frac{d^2}{N} + \frac{d^2}{B_y^2 N}\right)$ by the above choice of β .

We can thus apply Theorem 4 in the unregularized case ($\lambda = 0$) to obtain that, there exists a transformer with $\max_{i \in [L]} M^{(i)} \leq 3 \kappa \beta (4R + 4) = \max_{i \in [L]} (M^{(i)} + 4R + 4)$ (with $R = \max\{\beta B_x B_w, \beta B_y\}$), and number of layers

$$L = O\left(\log \frac{B_x B_w}{\beta} + \log \frac{B_x}{\beta} + \log \frac{B_y^2}{\beta}\right);$$

such that on $E_{\text{cov}} \setminus E_w$ (so that D is well-conditioned), we have (choosing the clipping radius in $\text{read}_y(\cdot) = \text{clip}_{B_y}(\text{read}_y(\cdot))$ to be B_y):

$$\text{read}_y(\text{TF}^0(\mathbf{H})) = \text{clip}_{B_y}(\mathbf{w}_{\text{LS}}^T \mathbf{x}_{N+1}) = \text{read}_y(\text{TF}^0(\mathbf{H})) \cdot \mathbf{w}_{\text{LS}}^T \mathbf{x}_{N+1} \quad \gamma = \frac{d^2}{N}; \quad (32)$$

We now bound the excess risk of the above transformer. Combining Proposition F.1 and (32), we have

$$E \left[\text{read}_y(\text{TF}^0(\mathbf{H})) \right] \leq \frac{1}{N_{\min}} \epsilon$$

$$\begin{aligned}
&= \mathbb{E} \left[\text{rad}_y(\text{TF}^0(\mathbf{H})) \right] y_{N+1}^2 \mathbb{1}_{E_{\text{cov}} \setminus E_w} + \mathbb{E} \left[\text{rad}_y(\text{TF}^0(\mathbf{H})) \right] y_{N+1}^2 \mathbb{1}_{(E_{\text{cov}} \setminus E_w)^c} \\
&\quad + 2 \mathbb{E} \left[\text{rad}_y(\text{TF}^0(\mathbf{H})) \right] \text{clip}_{B_y}(\mathbf{h}_{\text{LS}}; \mathbf{x}_{N+1})^2 \mathbb{1}_{E_{\text{cov}} \setminus E_w} \\
&\quad + 2 \mathbb{E} \left[\text{clip}_{B_y}(\mathbf{h}_{\text{LS}}; \mathbf{x}_{N+1}) \right] y_{N+1}^2 \mathbb{1}_{E_{\text{cov}} \setminus E_w} + 2B_y^2 = 10 \\
&\stackrel{(i)}{=} 2 \frac{d^2}{N} + L_P(\mathbf{w}_P^2) + O(B_y^2) + \frac{d^2}{N} + O(B_y^2) \\
&\quad L_P(\mathbf{w}_P^2) + O(B_y^2) + \frac{d^2}{N} = O\left(\frac{d^2}{N}\right) :
\end{aligned}$$

Above, (i) uses the approximation guarantee (32) as well as Proposition F.1(a) (with clipping radius B_y). This proves the desired excess risk guarantee.

Finally, under the canonical choice of parameters (23), the bounds for $L; M; \mathbb{J}; \mathbb{J}$ simplify to

$$L = O\left(\log \frac{N}{d}\right); \quad \max_{\ell \in [L]} M^{(\ell)} = 3; \quad \mathbb{J}; \mathbb{J} = O\left(\frac{d}{N}\right); \quad (33)$$

and the requirement for N simplifies to $N = O(d + \log(1/\delta)) = \Theta(d)$ (as $K = 1$). This proves the claim about the required N and L . \square

F.4 Proof of Corollary 6

Fix parameters $\delta, \epsilon > 0$ to be specified later and a large universal constant C_0 . Let us set

$$\begin{aligned}
&= \max_{0 \leq \ell \leq 2} \frac{d^{\ell+1}}{N^{\ell+1}}; \quad \epsilon = 25; \\
&B_w^2 := 1 + 2 \frac{\log(4/\delta)}{d}; \quad B_w = C_0(B_w^2 + \epsilon); \\
&B_x = C_0 \frac{d}{\log(N/\delta)}; \quad B_y = C_0(B_w^2 + \epsilon) \frac{d}{\log(N/\delta)};
\end{aligned}$$

Consider the following good events (below $\mathbf{y} = [y_i]_{i \in [M]} \in \mathbb{R}^N$ is given by $y_i = \mathbf{h}_{\text{LS}}(\mathbf{x}_i)$)

$$\begin{aligned}
E &= \|\mathbf{w}_{\text{ridge}}\|_2 \leq B_w; \quad \|\mathbf{x}_i\|_2 \leq \frac{d}{N}; \\
E_w &= \min(\mathbf{X}^T \mathbf{X} = N) \leq \max(\mathbf{X}^T \mathbf{X} = N) \leq \epsilon; \\
E_b &= \|\mathbf{x}_i\|_2 \leq B_x; \quad \|y_i\| \leq B_y; \\
E_{b;N+1} &= \|\mathbf{x}_{N+1}\|_2 \leq B_x; \quad \|y_{N+1}\| \leq B_y;
\end{aligned}$$

and we define $E := E \setminus E_w \setminus E_b \setminus E_{b;N+1}$. Under the event E , the problem (ICRidge) is well-conditioned and $\|\mathbf{w}_{\text{ridge}}\|_2 \leq B_w$ (by Lemma F.1).

Therefore, Theorem 4 implies that for $\epsilon = \frac{\delta}{2}$, there exists a $L = d \log(B_w/\epsilon) + 1$ -layer transformer with prediction $\hat{y}_{N+1} := \text{rad}_y(\text{TF}^0(\mathbf{H}))$ (clipped by B_y), such that under the good event E , we have $\hat{y}_{N+1} = \text{clip}_{B_y}(\mathbf{h}_{\text{LS}}; \mathbf{x}_{N+1})$ and $\|\mathbf{w}_{\text{ridge}}\|_2 \leq B_w$.

In the following, we show that \hat{y}_{N+1} is indeed the desired transformer (when δ and ϵ is suitably chosen). Notice that we have

$$\mathbb{E}(\hat{y}_{N+1} - y_{N+1})^2 = \mathbb{E} \mathbb{1}_{E^c} (\hat{y}_{N+1} - y_{N+1})^2 + \mathbb{E} \mathbb{1}_E (\hat{y}_{N+1} - y_{N+1})^2;$$

and we analyze these two parts separately.

Prediction risk under good event E . We first note that

$$\begin{aligned}
\mathbb{E} \mathbb{1}_E (\hat{y}_{N+1} - y_{N+1})^2 &= \mathbb{E} \mathbb{1}_E (\text{clip}_{B_y}(\mathbf{h}_{\text{LS}}; \mathbf{x}_{N+1}) - y_{N+1})^2 \\
&= \mathbb{E} \mathbb{1}_E (\mathbf{h}_{\text{LS}}(\mathbf{x}_{N+1}) - y_{N+1})^2;
\end{aligned}$$

where the inequality is because $y_{N+1} \in [B_y; B_y]$ under the good event E . Notice that by our construction, under the good event E , $\mathbf{w} = \mathbf{w}(D)$ depends only on the dataset D^7 . Therefore, we have $k\mathbf{w}(D) - \mathbf{w}_{\text{ridge}}(D)k \leq \epsilon$ as long as the event $E_0 := E \setminus E_w \setminus E_b$ holds for $(\mathbf{w}; D)$. Thus, under E_0 ,

$$\begin{aligned} \mathbb{E} \mathbb{1}_{E_0} \mathbb{E} g(\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}}(y_{N+1})^2 - \mathbf{w}; D) &= \mathbb{E} \mathbb{1}_{E_0} \mathbb{E} g(\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}(D)}(y_{N+1})^2 - \mathbf{w}; D) \\ &\quad + \mathbb{E} \mathbb{1}_{E_0} (\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}(D)}(y_{N+1})^2 - \mathbf{w}; D) \\ &= \mathbb{E} (\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}(D)}(y_{N+1})^2 - \mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}}(y_{N+1})^2) - \mathbf{w}; D + \epsilon^2 \\ &= k\mathbf{w}(D) - \mathbf{w}_{\text{ridge}}(D)k^2 + \epsilon^2; \end{aligned}$$

and we also have

$$\begin{aligned} k\mathbf{w}(D) - \mathbf{w}_{\text{ridge}}(D)k^2 &\leq \mathbf{w}_{\text{ridge}}^2 + 2\epsilon \mathbf{w}_{\text{ridge}} + \epsilon^2 \\ &\leq \mathbf{w}_{\text{ridge}}^2 + 2\epsilon \mathbf{w}_{\text{ridge}} + \epsilon^2; \end{aligned}$$

Recall that $2\text{BayesRisk} = \mathbb{E}_{\mathbf{w}; D} k\mathbf{w}_{\text{ridge}} - \mathbf{w}_{\text{ridge}}(D)k^2 + \epsilon^2$. Note that $2\text{BayesRisk} \leq 1 + \epsilon^2$ by definition. Therefore, we can conclude that

$$\mathbb{E} \mathbb{1}_{E_0} \mathbb{E} g(\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}}(y_{N+1})^2 - \mathbf{w}; D) \leq 2\text{BayesRisk} + 2\epsilon \mathbf{w}_{\text{ridge}} + \epsilon^2.$$

Prediction risk under bad event E^c . Notice that

$$\mathbb{E} \mathbb{1}_{E^c} \mathbb{E} g(\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}}(y_{N+1})^2 - \mathbf{w}; D) \leq \mathbb{P}(E^c) \mathbb{E}[(\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}}(y_{N+1})^2 - \mathbf{w}; D)^2].$$

We can upper bound $\mathbb{P}(E^c) = \mathbb{P}(E^c \cap [E_w^c \cap E_b^c \cap E_{b; N+1}^c])$ by Lemma B.1, Lemma B.2 and the sub-Gaussian tail bound:

$$\mathbb{P}(E^c) \leq \frac{1}{2} + \exp(-N/8); \quad \mathbb{P}(E_w^c) \leq 2 \exp(-N/8); \quad \mathbb{P}(E_b^c \cap E_{b; N+1}^c) \leq \frac{1}{4}.$$

Thus, as long as $N \geq 8 \log(12/\epsilon)$, we have $\mathbb{P}(E^c) \leq \epsilon$. Further, a simple calculation yields

$$\mathbb{E}(\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}}(y_{N+1})^2 - \mathbf{w}; D)^2 \leq 8\mathbb{E}y_{N+1}^4 + 8\mathbb{E}y_{N+1}^2 + 8B_y^2 + 8\mathbb{E}y_{N+1}^4.$$

Notice that $y_{N+1} | \mathbf{w} \sim \mathcal{N}(0; k\mathbf{w}_{\text{ridge}} - \mathbf{w}_{\text{ridge}}(D)k^2 + \epsilon^2)$, hence $\mathbb{E}y_{N+1}^4 = 3\mathbb{E}(k\mathbf{w}_{\text{ridge}} - \mathbf{w}_{\text{ridge}}(D)k^2 + \epsilon^2)^2 = 3(3 + 2\epsilon^2 + \epsilon^4) B_y^4$. Thus, we can conclude that

$$\mathbb{E} \mathbb{1}_{E^c} \mathbb{E} g(\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}}(y_{N+1})^2 - \mathbf{w}; D) \leq 4\epsilon B_y.$$

Choosing ϵ and δ . Combining the inequalities above, we have

$$\mathbb{E}(\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}}(y_{N+1})^2 - \mathbf{w}; D) \leq 2\text{BayesRisk} + 2\epsilon \mathbf{w}_{\text{ridge}} + \epsilon^2 + 4\epsilon B_y.$$

To ensure $\frac{1}{2} \mathbb{E}(\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}}(y_{N+1})^2 - \mathbf{w}; D) \leq \text{BayesRisk} + \epsilon$, we only need to take (ϵ, δ) so that the following constraints are satisfied:

$$\epsilon = \frac{1}{2} \min\{\delta, \delta^{\frac{1}{2}}\}; \quad 4\epsilon B_y \leq \frac{\delta}{2}; \quad N \geq 8 \log(12/\epsilon).$$

Therefore, it suffices to take $\delta = \frac{c_0}{\log^2(N)} \frac{\epsilon^2}{1 + \epsilon^2}$ for some small constant c_0 , then as long as

$$N \geq C \log \frac{2 + 1}{\epsilon} + C$$

our choice of ϵ and δ is feasible. Note that $\delta = O(1 + \epsilon^{-2})$, and hence under such choice of (ϵ, δ) , we have $L = O(\log(1/\delta))$ and $\mathbb{E} \mathbb{1}_{E^c} \mathbb{E} g(\mathbf{h}_{\mathbf{x}_{N+1}; \mathbf{w}}(y_{N+1})^2 - \mathbf{w}; D) = O(\delta)$. This is the desired result. \square

⁷We need this, as on E^c , the transformer output at this location could in principle depend additionally on \mathbf{x}_{N+1} , as (15) may not hold due to the potential unboundedness of its input. A similar fact will also appear in later proofs (for generalized linear models and Lasso).

Lemma F.1. Under the event $E \setminus E_w$, we have $\mathbf{w}_{\text{ridge}} = O(B_w^2 + \epsilon)$.

Proof of Lemma F.1. By the definition of $\mathbf{w}_{\text{ridge}}$ and recall that $d = d^2 = N$, we have $\mathbf{w}_{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + d^2 \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{y}$.

Therefore, we only need to prove the following fact: for any $\epsilon > 0$ and $\mathbf{b} = (\mathbf{X}^\top \mathbf{X} + d \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{y}$, we have

$$\|\mathbf{b} - \mathbf{w}_{\text{ridge}}\|_2 \leq B_w^2 + 10\epsilon \quad (34)$$

We now prove (34). Note that we have

$$\|\mathbf{b} - \mathbf{w}_{\text{ridge}}\|_2 = \|(\mathbf{X}^\top \mathbf{X} + d \mathbf{I}_d)^{-1} \mathbf{X}^\top (\mathbf{X} \mathbf{w}_{\text{ridge}} + \mathbf{y}) - \mathbf{b}\|_2 = \|\mathbf{B}_1 \mathbf{w}_{\text{ridge}} - \mathbf{b}\|_2 + \|\mathbf{B}_2 \mathbf{w}_{\text{ridge}}\|_2$$

where $\mathbf{B}_1 = \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + d \mathbf{I}_d)^{-1}$, $\mathbf{B}_2 = (\mathbf{X}^\top \mathbf{X} + d \mathbf{I}_d)^{-1} \mathbf{X}^\top$. Note that $\|\mathbf{B}_1\|_{\text{op}} \leq 1$ clearly holds, and under E we also have $\|\mathbf{B}_2\|_{\text{op}} \leq \frac{1}{\sqrt{N}}$. Therefore, it remains to bound the term $\|\mathbf{B}_2 \mathbf{w}_{\text{ridge}}\|_2$.

Consider the SVD decomposition of $\mathbf{X} = U V$, $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_d)$, and $U \in \mathbb{R}^{N \times d}$; $V \in \mathbb{R}^{d \times d}$ are orthonormal matrices. Then $\mathbf{B}_2 = V^\top (\Sigma^2 + d \mathbf{I}_d)^{-1} U^\top$, and hence

$$\|\mathbf{B}_2\|_{\text{op}} = \max_i \frac{\sigma_i}{\sigma_i^2 + d}$$

When $N \geq 36d$, we directly have $\|\mathbf{B}_2\|_{\text{op}} \leq \frac{1}{\sqrt{N}}$. Otherwise, we have $N < 36d$, and then for each $i \in [d]$, $\frac{\sigma_i}{\sigma_i^2 + d} \leq \frac{1}{\sqrt{N}}$. Hence, in this case we also have $\|\mathbf{B}_2\|_{\text{op}} \leq \frac{1}{\sqrt{N}}$. Combining the both cases completes the proof of (34). \square

G In-context learning of generalized linear models

As a natural generalization of linear regression, we now show that transformers can recover learn generalized linear models (GLMs) [53] (which includes logistic regression for linear classification as an important special case), by implementing the corresponding convex risk minimization algorithm in context, and achieve near-optimal excess risk under standard statistical assumptions.

Let $g: \mathbb{R} \rightarrow \mathbb{R}$ be a link function that is non-decreasing and C^2 -smooth. We consider the following convex empirical risk minimization (ERM) problem

$$\mathbf{w}_{\text{GLM}} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} \mathbb{E}_N(\mathbf{w}) := \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{x}_i; \mathbf{w}; y_i); \quad (\text{ICGLM})$$

where $\ell(t; y) := y t + \int_0^t g(s) ds$ is the convex (integral) loss associated with g . A canonical example of (ICGLM) is logistic regression, in which $g(t) = \text{logit}(t) := (1 + e^{-t})^{-1}$ is the sigmoid function, and the resulting $\ell(t; y) = \text{logit}(t; y) = y t + \log(1 + e^t)$ is the logistic loss.

The following result (proof in Appendix G.1) shows that, as long as the empirical risk \mathbb{E}_N satisfies strong convexity and bounded solution conditions (similar as in Theorem 4), transformers can approximately implement the ERM predictor $g(\mathbf{x}_{N+1}; \mathbf{w}_{\text{GLM}})$, with \mathbf{w}_{GLM} given by (ICGLM).

Theorem G.1 (Implementing convex risk minimization for GLMs). For any $0 < \epsilon < 1$ with $\epsilon := \frac{1}{2}$, $B_w > 0$; $B_x > 0$, $\epsilon := L_g B_x^2 = \epsilon + 1$ and $\epsilon < B_w = 2$, there exists an attention-only transformer TF^ϵ with

$$L = d \log(L_g B_w B_x) \epsilon + 1; \quad \max_{\mathbf{w} \in [L]} M(\mathbf{w}) \leq C_g \frac{2}{w} \epsilon^{-2}; \quad \|\mathbf{w}_{\text{GLM}}\|_2 \leq R + \epsilon^{-1} C_g;$$

(where $L_g := \sup_t |g'(t)|$, $R := \max\{B_x B_w, B_y\}$, and $C_g > 0$ is a constant that depends only on R and the C^2 -smoothness of g within $[-R; R]$), such that the following holds. On any input data $(D; \mathbf{x}_{N+1})$ such that

$$\min(\epsilon^{-2} \mathbb{E}_N(\mathbf{w})) \leq \max(\epsilon^{-2} \mathbb{E}_N(\mathbf{w})) \leq \epsilon^{-2} \quad \text{for all } \mathbf{w} \in \mathcal{B}_2(B_w); \quad \|\mathbf{w}_{\text{GLM}}\|_2 \leq B_w = 2; \quad (35)$$

$\text{TF}^0(\mathbf{H}^{(0)})$ approximately implements (ICGLM): We have $\mathbf{h}_{N+1}^{(L+1)} := [\mathbf{x}_{N+1}; \mathbf{y}_{N+1}; \mathbf{w}; 1; 1]$, where

$$\mathbf{y}_{N+1} = g(\mathbf{h}_{N+1}; \mathbf{w}_{\text{GLM}}) \quad \text{“”}$$

In Theorem G.1, the number of heads scales as $\Theta(1/\epsilon^2)$ as opposed to (1) as in ridge regression (Theorem 4), due to the fact that the gradient of the loss is in general a smooth function that can be only approximately expressed as a sum-of-relus (cf. Definition D.1 & Lemma B.5) rather than exactly expressed as in the case for the square loss.

In-context prediction power We next show that (proof in Appendix G.2) the transformer constructed in Theorem G.1 achieves desirable statistical power if the in-context data distribution satisfies standard statistical assumptions for learning GLMs. Let $L_P(\mathbf{w}) := \mathbb{E}_{(\mathbf{x}; y) \sim P} [(\mathbf{h}\mathbf{w}; \mathbf{x}; y)]$ denote the corresponding population risk for any distribution P of $(\mathbf{x}; y)$. When P is *realizable* by a generalized linear model of link function g and parameter \mathbf{w} in the sense that $\mathbb{E}_P[y|\mathbf{x}] = g(\mathbf{h}; \mathbf{x}; \mathbf{w})$, it is a standard result that \mathbf{w} is indeed a minimizer of L_P [42] (see also [6, Appendix A.3]).

Theorem G.2 (Statistical guarantee for generalized linear models). *For any fixed set of parameters defined in Assumption B, there exists a transformer with $L = O(\log(N))$ layers and $\max_{i \in [L]} M^{(i)} = O(d^\beta N)$, such that for any distribution P satisfying Assumption B with those parameters, as long as $N = O(d)$, that outputs $\mathbf{y}_{N+1} = \text{read}_y(\text{TF}(\mathbf{H}))$ and $\mathbf{w} = \text{read}_w(\text{TF}(\mathbf{H})) \in \mathbb{R}^d$ (for another read-out function read_w) satisfying the following.*

(a) \mathbf{w} achieves small excess risk under the population loss, i.e. for the linear prediction $\mathbf{y}_{N+1}^{\text{lin}} := \mathbf{h}_{N+1}; \mathbf{w}$,

$$\mathbb{E}_{(\mathcal{D}; \mathbf{x}_{N+1}; \mathbf{y}_{N+1}) \sim P} (\mathbf{y}_{N+1}^{\text{lin}} - \mathbf{y}_{N+1})^2 = \min_{\mathbf{w}} L_P(\mathbf{w}) - O(d=N) \quad (36)$$

(b) (Realizable setting) If there exists a $\mathbf{w} \in \mathbb{R}^d$ such that under P , $\mathbb{E}[y|\mathbf{x}] = g(\mathbf{h}; \mathbf{x}; \mathbf{w})$ almost surely, then

$$\mathbb{E}_{(\mathcal{D}; \mathbf{x}_{N+1}; \mathbf{y}_{N+1}) \sim P} (\mathbf{y}_{N+1} - \mathbf{y}_{N+1})^2 = \mathbb{E}_{(\mathbf{x}_{N+1}; \mathbf{y}_{N+1}) \sim P} (g(\mathbf{h}; \mathbf{x}_{N+1}; \mathbf{w}) - \mathbf{y}_{N+1})^2 + O(d=N); \quad (37)$$

or equivalently, $\mathbb{E}[(\mathbf{y}_{N+1} - \mathbb{E}[\mathbf{y}_{N+1}|\mathbf{x}_{N+1}])^2] = O(d=N)$.

Above, $O(\cdot)$ hides constants that depend polynomially on the parameters in Assumption B. Similar as in Corollary 5, the $O(d=N)$ excess risk obtained here matches the optimal (fast) rate for typical learning problems with d parameters and N samples [87].

Assumption B (Well-posedness for learning GLMs). *We assume that there is some $B > 0$ such that for any $t \geq [B; B]$, $g'(t) = g > 0$.*

We also assume that for each $i \geq [N + 1]$, $(\mathbf{x}_i; y_i)$ is independently sampled from P such that the following holds.

(a) Under the law $(\mathbf{x}; y) \sim P$, We have $\mathbf{x} \sim \text{SG}(K_x)$, $y \sim \text{SG}(K_y)$ and $g(\mathbf{h}\mathbf{w}; \mathbf{x}; y) \sim \text{SG}(K_y)$ $\forall \mathbf{w} \in \mathbb{B}_2(B_w)$.

(b) For some $\epsilon > 0$, it holds that

$$\mathbb{E}[\mathbf{1}^T \mathbf{f}(\mathbf{x}) \mathbf{w} \mathbf{j} - B = 2g\mathbf{x}\mathbf{x}^T] \leq \epsilon \mathbf{1}_d \quad \forall \mathbf{w} \in \mathbb{B}_2(B_w);$$

(c) For $\mathbf{w}^* = \arg \min_{\mathbf{w}} L_P$, it holds $k_2(\mathbf{w}^*) \leq B_w = 4$.

Applying Theorem G.2 to logistic regression, we have the following result as a direct corollary. Below, the Gaussian input assumption is for convenience only and can be generalized to e.g. sub-Gaussian input.

Corollary G.1 (In-context logistic regression). *Consider any in-context data distribution P satisfying*

$$\mathbf{x} \sim \mathcal{N}(0; \mathbf{I}_d); \quad y \geq \tau; 1-g; \quad \arg \min_{\mathbf{w} \in \mathbb{R}^d} L_P(\mathbf{w}) \in \mathbb{B}_2(B_w^*);$$

For the link function $g = \text{logit}$ and $B_w^ = O(1)$, we can choose $B_w; B; \epsilon; g; L_g; \epsilon; K_x; K_y = O(1)$ so that Assumption B holds. In that case, when $N = O(d)$, there exists a transformer with $L = O(\log(N))$ layers, such that for any P considered above,*

(a) The estimation $\hat{\mathbf{w}} = \hat{\text{rad}}_w(\text{TF}(\mathbf{H}))$ outputted by $\hat{\text{rad}}_w$ achieves excess risk bound (36).

(b) (Realizable setting) Consider the logistic in-context data distribution

$$P^{\text{log}}: \quad \mathbf{x} \sim N(\mathbf{0}; \mathbf{I}_d); \quad y|\mathbf{x} \sim \text{Bernoulli}(g(h(\mathbf{x}; \mathbf{w})))$$

Then, for any distribution $P = P^{\text{log}}$ with $k_2 \leq B_w^2$, the prediction $\hat{y}_{N+1} = \hat{\text{rad}}_y(\text{TF}(\mathbf{H}))$ of $\hat{\text{rad}}_y$ additionally achieves the square loss excess risk (37).

G.1 Proof of Theorem G.1

Let us fix parameters $\epsilon_g > 0$ and $T > 0$ (that we specify later in proof).

Define $R = \max\{B_x B_w, B_y\} \epsilon_g$ and

$$C_g := \max_{i=0,1,2} R^i \max_{s \in [-B; B]} g^{(i)}(s) \quad ;$$

By Proposition B.1, g is $(\epsilon_g; M; R; C)$ with

$$C = O(C_g); \quad M = O(C_g^2 \epsilon_g^{-2} \log(1 + C_g \epsilon_g^{-1})) \quad ;$$

Therefore, we can invoke Theorem D.1 to obtain that, as long as $2T \epsilon_g \leq B_w$, there exists a T -layer attention-only transformer $\hat{\text{rad}}_y^{(1:T)}$ with M heads per layer, such that for any input \mathbf{H} of format (3) and satisfies (35), its last layer outputs $\mathbf{h}_i^{(T)} = [\mathbf{x}_i; y_i; \hat{\mathbf{w}}^T; \mathbf{0}_{D-2d-3}; 1; t_i]$, such that

$$\hat{\mathbf{w}}^T = \mathbf{w}_{\text{GD}}^T \quad \epsilon_g \quad (L^{-1} B_x);$$

where $\hat{\mathbf{w}}_{\text{GD}} = \hat{\mathbf{w}}_{\text{GD}}^{g \in [L]}$ is the sequence of gradient descent iterates with stepsize ϵ_g^{-1} and initialization $\mathbf{w}_{\text{GD}}^0 = \mathbf{0}$. Notice that Proposition B.2 implies (with $\epsilon_g = \epsilon_g$)

$$\|\mathbf{w}_{\text{GD}}^T - \mathbf{w}_{\text{GLM}}^T\|_2 \leq \exp(-T \epsilon_g) k_{\text{GLM}} \exp(-T \epsilon_g) \frac{B_w}{2} := \epsilon_o;$$

Furthermore, we can show that (similar to the proof of Theorem D.1 (b)), there exists a single attention layer $\hat{\text{rad}}_y^{(T+1)}$ with M heads such that it outputs $\mathbf{h}_{N+1}^{(T+1)} = [\mathbf{x}_{N+1}; y_{N+1}; \hat{\mathbf{w}}^T; \mathbf{0}_{D-2d-3}; 1; 0]$, where $y_{N+1} = g(\mathbf{x}_{N+1}; \hat{\mathbf{w}}^T) \pm \epsilon_g$.

In the following, we show that for suitably chosen $(T; \epsilon_g)$, $\hat{\text{rad}}_y^{(1:T); \hat{\text{rad}}_y^{(T+1)}}$ is the desired transformer. First notice that its output $\mathbf{h}_{N+1}^{(T+1)} = [\mathbf{x}_{N+1}; y_{N+1}; \hat{\mathbf{w}}^T; \mathbf{0}_{D-2d-3}; 1; 0]$ satisfies

$$\begin{aligned} |y_{N+1} - g(\mathbf{x}_{N+1}; \mathbf{w}_{\text{GLM}}^T)| &\leq |y_{N+1} - g(\mathbf{x}_{N+1}; \hat{\mathbf{w}}^T)| + L_g \|\mathbf{x}_{N+1}; \hat{\mathbf{w}}^T - \mathbf{w}_{\text{GLM}}^T\|_2 \\ &\leq \epsilon_g + L_g B_x \|\hat{\mathbf{w}}^T - \mathbf{w}_{\text{GD}}^T\|_2 + L_g B_x \|\mathbf{w}_{\text{GD}}^T - \mathbf{w}_{\text{GLM}}^T\|_2 \\ &\leq \epsilon_g (1 + L_g B_x T^{-1} B_x) + L_g B_x \epsilon_o; \end{aligned}$$

Therefore, for any fixed $\epsilon_g > 0$, we can take

$$T = d \log(L_g B_x B_w \epsilon_g^{-1}); \quad \epsilon_g = \frac{1}{2} \frac{\epsilon_o}{1 + T (L_g B_x^2 \epsilon_g^{-1})};$$

so that the we construct above ensures $|y_{N+1} - g(\mathbf{x}_{N+1}; \mathbf{w}_{\text{GLM}}^T)| \leq \epsilon_g$ for any input \mathbf{H} that satisfies (35). The upper bound on ϵ_g follows immediately from Theorem D.1. \square

G.2 Proof of Theorem G.2

We summarize some basic and useful facts about GLM in the following theorem. Its proof is presented in Appendix G.3 - G.6.

Theorem G.3. Under Assumption B, the following statements hold with universal constant C_0 and constant $C_1; C_2$ that depend only on the parameters $(K_x; K_y; B; B_w; \epsilon_x; L_g; \epsilon_g)$.

(a) As long as $N \geq C_1 d$, the following event happens with probability at least $1 - 2e^{-N=C_1}$:

$$E_w: \frac{1}{8} g \leq \min(r^2 \mathbb{E}_N(\mathbf{w})) \leq \max(r^2 \mathbb{E}_N(\mathbf{w})) \leq 8L_g K_x^2; \quad \forall \mathbf{w} \in \mathcal{B}_2(B_w)$$

(b) For any $\epsilon > 0$, we have with probability at least $1 - \epsilon$ that

$$\|\mathbf{w}_{\text{stat}} - \mathbb{E}_N(\mathbf{w})\|_2 \leq \sup_{\mathbf{w} \in \mathcal{B}_2(B_w)} \|\mathbf{r}_w \mathbb{E}_N(\mathbf{w}) - \mathbf{r}_w \mathbb{E}[\mathbb{E}_N(\mathbf{w})]\|_2 \leq C_0 K_x K_y \max\left(\frac{r}{N}, \frac{d + \log(1/\epsilon)}{N}\right);$$

where we denote $B_w = \log(2 + L_g K_x^2 B_w) = K_y$.

(c) Condition on (a) holds and $N \geq C_2 d$, the event $E_r := \|\mathbf{w}_{\text{GLM}}\|_2 \leq B_w = 2g$ happens with probability at least $1 - e^{-N=C_2}$.

(d) For any $\mathbf{w} \in \mathcal{B}_2(B_w)$, it holds that

$$L_p(\mathbf{w}) - L_p(\mathbf{w}_{\text{stat}}) \leq \frac{4}{g} \|\mathbf{w}_{\text{stat}} - \mathbf{r} \mathbb{E}_N(\mathbf{w})\|_2^2;$$

(e) (Realizable setting) As long as $\mathbf{w}_{\text{GLM}} \in \mathcal{B}_2(B_w)$, it holds that

$$\mathbb{E}_x(g(\mathbf{h}\mathbf{x}; \mathbf{w}_{\text{GLM}}) - g(\mathbf{h}\mathbf{x}; \mathbf{w}_{\text{stat}}))^2 \leq \frac{L_g}{x} \|\mathbf{w}_{\text{stat}}\|_2^2;$$

Therefore, we can set

$$B_x = \frac{g}{8} \sqrt{\frac{1}{d \log(N)}}; \quad B_y = 8L_g K_x^2 \sqrt{\frac{1}{d \log(N)}};$$

Consider the following good events

$$\begin{aligned} E_b &= \|\mathbf{x}_i\|_2 \leq B_x; \quad \forall i \in [N]; \\ E_{b;N+1} &= \|\mathbf{x}_{N+1}\|_2 \leq B_x; \\ E &= E_r \setminus E_w \setminus E_b \setminus E_{b;N+1}; \end{aligned}$$

Under the event E and our choice of B_x, B_y , the problem (ICGLM) is well-conditioned (i.e. (35) holds).

Theorem G.1 implies that there exists a transformer \mathcal{T} such that for any input \mathbf{H} of the form (3), \mathcal{T} outputs $\mathbf{h}'_{N+1} = [\mathbf{x}_{N+1}; \mathbf{y}_{N+1}; \mathbf{w}; \mathbf{0}_{D-2d-3}; 1; 0]$, such that the output is given by $\mathbf{y}_{N+1} = \text{read}_y(\mathcal{T}(\mathbf{H})) = \text{clip}_{B_y}(\mathbf{y}_{N+1})$ and $\mathbf{w} = \text{read}_w(\mathcal{T}(\mathbf{H})) := \text{Proj}_{\mathcal{B}_2(B_w)}(\mathbf{w})$, and the following holds on the good event E :

- (a) $\mathbf{y}_{N+1} = f_D(\mathbf{x}_{N+1})$, where $f_D = A(D)$ is a predictor such that $\|f_D(\mathbf{x}) - g(\mathbf{h}\mathbf{x}; \mathbf{w}_{\text{GLM}})\|_2 \leq \epsilon$ for all $\mathbf{x} \in \mathcal{B}_2(B_x)$.
- (b) $\mathbf{w} = \mathbf{w}(D) \in \mathcal{B}_2(B_w)$ depends only on D (by the proof of Theorem G.1 and Theorem D.1), such that $\|\mathbf{r} \mathbb{E}_N(\mathbf{w}) - \mathbf{w}\|_2 \leq \frac{\epsilon}{L_g B_w}$.

In the following, we show that \mathcal{T} constructed above fulfills both (a) & (b) of Theorem G.2. The bounds on number of layers and heads and $\mathcal{J} \leq \mathcal{J}$ follows from plugging our choice of B_x, B_y in our proof of Theorem G.1.

Proof of Theorem G.2 (a). Notice that under the good event E , we have $\mathbf{w} = \mathbf{w} = \mathbf{w}(D)$ depends only on D . Then we have

$$\begin{aligned} & \mathbb{E}_{(D; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})} \|\mathbf{y}_{N+1}^{\text{lin}} - \mathbf{y}_{N+1}\|_2^2 \\ &= \mathbb{E}_{(D; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})} \|\mathbf{f}_D(\mathbf{x}_{N+1}) - \mathbf{y}_{N+1}\|_2^2 + \mathbb{E}_{(D; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})} \|\mathbf{f}_D(\mathbf{x}_{N+1}) - \mathbf{g}(\mathbf{h}\mathbf{x}_{N+1}; \mathbf{w}(D))\|_2^2 \\ &= \mathbb{E}_{(D; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})} \|\mathbf{f}_D(\mathbf{x}_{N+1}) - \mathbf{g}(\mathbf{h}\mathbf{x}_{N+1}; \mathbf{w}(D))\|_2^2 + \mathbb{E}_{(D; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})} \|\mathbf{f}_D(\mathbf{x}_{N+1}) - \mathbf{y}_{N+1}\|_2^2; \end{aligned}$$

Thus, we can consider $E_0 = E_r \setminus E_w \setminus E_b$, and then

$$\mathbb{E}_{(D; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})} \|\mathbf{f}_D(\mathbf{x}_{N+1}) - \mathbf{g}(\mathbf{h}\mathbf{x}_{N+1}; \mathbf{w}(D))\|_2^2$$

$$= E_{(\mathcal{D}; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})}[1fE_0g(\mathbf{h}\mathbf{x}_{N+1}; \mathbf{w}(D); i; \mathbf{y}_{N+1})] E_{(\mathcal{D}; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})}[1fE_0 Eg(\mathbf{h}\mathbf{x}_{N+1}; \mathbf{w}(D); i; \mathbf{y}_{N+1})]$$

$$= E_{(\mathcal{D}; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})}[1fE_0gL_p(\mathbf{w}(D))] E_{(\mathcal{D}; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})}[1fE_0 Eg(\mathbf{h}\mathbf{x}_{N+1}; \mathbf{w}(D); i; \mathbf{y}_{N+1})];$$
 where the second equality follows from $L_p(\mathbf{w}(D)) = E_{(\mathbf{x}_{N+1}; \mathbf{y}_{N+1})|D}(\mathbf{h}\mathbf{x}_{N+1}; \mathbf{w}(D); i; \mathbf{y}_{N+1})$. Therefore,

$$\begin{aligned}
 & E_{(\mathcal{D}; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})}(\mathcal{Y}_{N+1}^{\text{lin}}; \mathbf{y}_{N+1}) E_{\mathcal{D}}[1fE_0gL_p(\mathbf{w}(D))] \\
 &= E_{(\mathcal{D}; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})} \frac{1fE^c g(\mathcal{Y}_{N+1}^{\text{lin}}; \mathbf{y}_{N+1}) E_{(\mathcal{D}; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})}[1fE_0 Eg(\mathbf{h}\mathbf{x}_{N+1}; \mathbf{w}(D); i; \mathbf{y}_{N+1})]}{2 P(E^c) \max E(\mathcal{Y}_{N+1}^{\text{lin}}; \mathbf{y}_{N+1})^4; E[(\mathbf{h}\mathbf{x}_{N+1}; \mathbf{w}(D); i; \mathbf{y}_{N+1})^4]} = O\left(\frac{B^2}{N^5}\right);
 \end{aligned}$$

where the last line follows from Cauchy inequality and the fact $P(E^c) = O(N^{-10})$, and B is defined in Lemma G.1.

Notice that by Theorem G.3 (d), we have

$$E_{\mathcal{D}}[1fE_0g(L_p(\mathbf{w}) \inf L_p)] \leq \frac{4}{g} E[\sigma_{\text{stat}}^2] + E[1fE_0g] r \mathbb{E}_N(\mathbf{w})^2;$$

and by Theorem G.3 (b) and taking integration over $\gamma > 0$, we have

$$E[\sigma_{\text{stat}}^2] \leq O(1) K_x^2 K_y^2 \frac{d}{N} + \frac{d}{N}^2;$$

Also, we have $\inf L_p = L_p(\gamma) \leq B$ by Lemma G.1. Therefore, we can conclude that

$$E_{(\mathcal{D}; \mathbf{x}_{N+1}; \mathbf{y}_{N+1})}(\mathcal{Y}_{N+1}^{\text{lin}}; \mathbf{y}_{N+1}) \leq \inf L_p + O(1) \left[\frac{K_x^2 K_y^2 d}{g} + \frac{K_x^4}{g} \sigma^2 + \frac{B^2}{N^5} \right];$$

Taking $\sigma^2 = \frac{K_y^2}{B_w K_x^2} \frac{d}{N}$ completes the proof. \square

Proof of Theorem G.2 (b). Similar to the proof of Corollary 6, we have

$$\begin{aligned}
 E(\mathcal{Y}_{N+1} \mathbf{y}_{N+1})^2 &= E[1fEg(\mathcal{Y}_{N+1} \mathbf{y}_{N+1})^2] + E[1fE^c g(\mathcal{Y}_{N+1} \mathbf{y}_{N+1})^2] \\
 &= E[1fEg(\mathcal{Y}_{N+1} \mathbf{y}_{N+1})^2] + \frac{P(E^c) E(\mathcal{Y}_{N+1} \mathbf{y}_{N+1})^4}{P(E^c)};
 \end{aligned}$$

where the inequality follows from $\mathcal{Y}_{N+1} \geq [B_y; B_y]$ on event E . For the first part, we have

$$\begin{aligned}
 E[1fEg(\mathcal{Y}_{N+1} \mathbf{y}_{N+1})^2] &= E[1fEg(f_{\mathcal{D}}(\mathbf{x}_{N+1}) \mathbf{y}_{N+1})^2] \\
 &= E_{\mathcal{D}}[1fE_0g E_{(\mathbf{x}; \mathbf{y}) \sim P}(\mathbf{f}k\mathbf{x}k_2 B_x g(f_{\mathcal{D}}(\mathbf{x}) \mathbf{y}))^2];
 \end{aligned}$$

where we use the fact that the conditional distribution of $(\mathbf{x}_{N+1}; \mathbf{y}_{N+1})|D$ agrees with P . Thus,

$$\begin{aligned}
 & E[1fEg(\mathcal{Y}_{N+1} \mathbf{y}_{N+1})^2] = E_{(\mathbf{x}; \mathbf{y}) \sim P}(g(\mathbf{h}\mathbf{x}; i) \mathbf{y})^2 \\
 &= E_{\mathcal{D}}[1fE_0g E_{(\mathbf{x}; \mathbf{y}) \sim P}(\mathbf{f}k\mathbf{x}k_2 B_x g(f_{\mathcal{D}}(\mathbf{x}) \mathbf{y}))^2] = E_{(\mathbf{x}; \mathbf{y}) \sim P}(g(\mathbf{h}\mathbf{x}; i) \mathbf{y})^2 \\
 &= E_{\mathcal{D}}[1fE_0g E_{\mathbf{x}}(\mathbf{f}k\mathbf{x}k_2 B_x g(f_{\mathcal{D}}(\mathbf{x}) g(\mathbf{h}\mathbf{x}; i)))^2] \\
 &= 2E_{\mathcal{D}}[1fE_0g E_{\mathbf{x}}(\mathbf{f}k\mathbf{x}k_2 B_x g(f_{\mathcal{D}}(\mathbf{x}) g(\mathbf{h}\mathbf{x}; \mathbf{w}_{\text{GLM}} i)))^2] + 2E_{\mathcal{D}}[1fE_0g E_{\mathbf{x}}(g(\mathbf{h}\mathbf{x}; \mathbf{w}_{\text{GLM}} i) g(\mathbf{h}\mathbf{x}; i))^2] \\
 &= 2\sigma^2 + \frac{2Lg}{x} E[\sigma_{\text{stat}}^2] = 2\sigma^2 + O(1) \frac{Lg K_x^2 K_y^2 d}{x g N}.
 \end{aligned}$$

For the second part, we know $P(E^c) = O(N^{-10})$ and

$$E(\mathcal{Y}_{N+1} \mathbf{y}_{N+1})^4 \leq 8E\mathcal{Y}_{N+1}^2 + 8E\mathbf{y}_{N+1}^4 = O(B_y^4);$$

In conclusion, we have

$$E(\mathcal{Y}_{N+1} \mathbf{y}_{N+1})^2 \leq E(\mathcal{Y}_{N+1} g(\mathbf{h}\mathbf{x}_{N+1}; i))^2 + 2\sigma^2 + O(1) \frac{Lg K_x^2 K_y^2 d}{x g N} + O\left(\frac{B_y^2}{N^5}\right);$$

Taking $\sigma^2 = \frac{Lg K_x^2 K_y^2 d}{x g N}$ completes the proof. \square

Lemma G.1. Suppose that $\mathbf{x} \in \text{SG}(K_x)$, $y \in \text{SG}(K_y)$, and \mathbf{w} is a (possibly random) vector such that $\|\mathbf{w}\|_2 \leq B_w$. Then

$$\mathbb{E} \left[\langle \mathbf{h}(\mathbf{x}; \mathbf{w}; y) \rangle^4 \right] \leq O(L_g K_x^2 B_w^2 d + K_x K_y B_w d) =: B:$$

Proof. Notice that by our assumption, $g(0) \leq 2K_y$. Therefore, by the definition of $\langle \cdot \rangle$,

$$\langle \mathbf{h}(\mathbf{x}; \mathbf{w}; y) \rangle = y t + \int_0^t g(s) ds = \int_0^t (g(s) - g(0)) ds + \int_0^t (g(0) + y) ds = \int_0^t (g(s) - g(0)) ds + t(2K_y + jy) + 2L_g t^2:$$

The proof is then done by bounding the moment by $\mathbb{E} |jy|^8 \leq O(K_y^8)$ and $\mathbb{E} \langle \mathbf{h}(\mathbf{x}; \mathbf{w}; y) \rangle^8 \leq B_w^8 \mathbb{E} \|\mathbf{x}\|_2^8 \leq O(d^4 B_w^8 K_x^8)$, which is standard (by utilizing the tail bound of sub-Gaussian/sub-Exponential random variable). \square

G.3 Proof of Theorem G.3 (a)

We begin with the upper bound on $\max(\mathbf{r}^2 \underline{\mathbf{p}}_N(\mathbf{w}))$. By Lemma B.3, as long as $N \geq C_0 d$, the following event

$$E_{w,0}: \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top \text{op} \leq 8K^2:$$

happens with probability at least $1 - \exp(-N/C_0)$. By the assumption that $\sup |g'| \leq L_g$, it is clear that when $E_{w,0}$ holds, we have $\max(\mathbf{r}^2 \underline{\mathbf{p}}_N(\mathbf{w})) \leq 8L_g K_x^2 \|\mathbf{w}\|_2 \leq 8L_g K_x^2 B_w$.

In the following, we analyze the quantity $\max(\mathbf{r}^2 \underline{\mathbf{p}}_N(\mathbf{w}))$. We have to invoke the following covering argument (see e.g. [85, Section 4.1.1]).

Lemma G.2. Suppose that V is a ϵ -covering of S^{d-1} with $\epsilon \in [0, 1)$. Then the following holds:

1. For any $d \times d$ symmetric matrix A , $\|A\|_{\text{op}} \leq \frac{1}{1-2\epsilon} \max_{\mathbf{v} \in V} \mathbf{v}^\top A \mathbf{v}$ and

$$\min(A) \leq \min_{\mathbf{v} \in V} \mathbf{v}^\top A \mathbf{v} \leq 2\epsilon \|A\|_{\text{op}}$$

2. For any vector $\mathbf{x} \in \mathbb{R}^d$, $\|\mathbf{x}\|_2 \leq \frac{1}{1-2\epsilon} \max_{\mathbf{v} \in V} \langle \mathbf{v}; \mathbf{x} \rangle$.

Notice that

$$\begin{aligned} \mathbf{r}^2 \underline{\mathbf{p}}_N(\mathbf{w}) &= \frac{1}{N} \sum_{i=1}^N g'(\langle \mathbf{h}(\mathbf{w}; \mathbf{x}_i) \rangle) \mathbf{x}_i \mathbf{x}_i^\top = \frac{1}{N} \sum_{i=1}^N g'(\langle \mathbf{h}(\mathbf{w}; \mathbf{x}_i) \rangle - B) \mathbf{x}_i \mathbf{x}_i^\top \\ &= \frac{1}{N} \sum_{i=1}^N g' \left(1 - \frac{\langle \mathbf{h}(\mathbf{w}; \mathbf{x}_i) \rangle}{B} \right) \mathbf{x}_i \mathbf{x}_i^\top. \end{aligned}$$

Therefore, we can define $h(t) := (B - \langle \mathbf{h}(\mathbf{w}; \mathbf{x}_i) \rangle)_+$ (which is a 1-Lipschitz function), and we have

$$\mathbf{r}^2 \underline{\mathbf{p}}_N(\mathbf{w}) = \frac{g}{B} \frac{1}{N} \sum_{i=1}^N h(\langle \mathbf{h}(\mathbf{w}; \mathbf{x}_i) \rangle) \mathbf{x}_i \mathbf{x}_i^\top =: \frac{1}{N} \sum_{i=1}^N \underbrace{h(\langle \mathbf{h}(\mathbf{w}; \mathbf{x}_i) \rangle)}_{=: A(\mathbf{w})} \mathbf{x}_i \mathbf{x}_i^\top$$

In the following, we pick a ϵ -covering V of S^{d-1} such that $|V| \leq (3/\epsilon)^d$ (we will specify ϵ later in proof). Then for any $\mathbf{w} \in B_2(B_w)$,

$$\min(A(\mathbf{w})) \leq \min_{\mathbf{v} \in V} \mathbf{v}^\top A(\mathbf{w}) \mathbf{v} \leq 2\epsilon \|A(\mathbf{w})\|_{\text{op}}$$

By our definition of $A(\mathbf{w})$, we have (for any fixed B_w)

$$\min_{\mathbf{v} \in V} \mathbf{v}^\top A(\mathbf{w}) \mathbf{v} = \min_{\mathbf{v} \in V} \frac{1}{N} \sum_{i=1}^N h(\langle \mathbf{h}(\mathbf{w}; \mathbf{x}_i) \rangle) \langle \mathbf{v}; \mathbf{x}_i \rangle^2$$

$$\min_{\mathbf{v} \in \mathcal{V}} \frac{1}{N} \sum_{i=1}^N h(\mathbf{w}; \mathbf{x}_i) \min_{\mathbf{v} \in \mathcal{V}} \sum_{i=1}^N h(\mathbf{w}; \mathbf{x}_i)^2; B_{xv}^2$$

$$\min_{\mathbf{v} \in \mathcal{V}} E[U_v(\mathbf{w})] + \min_{\mathbf{v} \in \mathcal{V}} (U_v(\mathbf{w}) - E[U_v(\mathbf{w})]):$$

By Lemma G.3, we can choose $B_{xv} = K_x(15 + \log(K_x^2 = x))$, and then $E[U_v(\mathbf{w})] \leq 3B_{xv} = 8$. Thus, combining the inequalities above, we can take $\nu = \frac{128K_x^2}{x}$ in the following, so that under event $E_{w,0}$,

$$\min_{\mathbf{v} \in \mathcal{V}} (r^2 \mathbb{E}_N(\mathbf{w})) \leq \frac{g}{8} + \frac{g}{B} \frac{B}{16} \max_{\mathbf{v} \in \mathcal{V}} (E[U_v(\mathbf{w})] - U_v(\mathbf{w})) :$$

In the following, we consider the random process $\bar{U}_v(\mathbf{w}) := U_v(\mathbf{w}) - E[U_v(\mathbf{w})]$, which is zero-mean and indexed by $\mathbf{w} \in \mathcal{B}_2(B_w)$. For any fixed \mathbf{v} , consider applying Proposition B.4 to the random process $\bar{U}_v(\mathbf{w})$. We need to verify the preconditions:

- (a) With norm $\langle \mathbf{w}; \mathbf{w}' \rangle = k\mathbf{w} \cdot \mathbf{w}'$, $\log N(\mathcal{B}(\mathbf{w}; r)) \leq d \log(2Ar)$ with constant $A = 2$;
- (b) Let $f(\mathbf{x}; \mathbf{w}) := h(\mathbf{w}; \mathbf{x}) \min_{\mathbf{v} \in \mathcal{V}} \sum_{i=1}^N h(\mathbf{w}; \mathbf{x}_i)^2; B_{xv}^2$, then $|f(\mathbf{x}; \mathbf{w})| \leq B_{xv}^2$ and hence in $\text{SG}(CB_{xv}^2)$ for any random \mathbf{x} ;
- (c) For $\mathbf{w}, \mathbf{w}' \in \mathcal{W}$, we have $|h(\mathbf{w}; \mathbf{x}) - h(\mathbf{w}'; \mathbf{x})| \leq |h(\mathbf{w}; \mathbf{x}) - h(\mathbf{w}'; \mathbf{x})|$. Hence, because $\mathbf{x} \in \text{SG}(K_x)$, the random variable $h(\mathbf{w}; \mathbf{x}) - h(\mathbf{w}'; \mathbf{x})$ is $\text{SG}(CK_x k\mathbf{w} \cdot \mathbf{w}' k_2)$, and the random variable $f(\mathbf{x}; \mathbf{w}) - f(\mathbf{x}; \mathbf{w}')$ is $\text{SG}(CK_x B_{xv}^2 k\mathbf{w} \cdot \mathbf{w}' k_2)$.

Therefore, we can apply Proposition B.4 to obtain that with probability $1 - \epsilon$, it holds

$$\sup_{\mathbf{w}} \bar{U}_v(\mathbf{w}) \leq C' B_{xv}^2 \frac{d \log(2g) + \log(1 - \epsilon)}{N};$$

where we denote $g = 1 + K_x B_w = B$. Setting $\epsilon = \nu/V$ and taking the union bound over $\mathbf{v} \in \mathcal{V}$, we obtain that with probability at least $1 - \epsilon$,

$$\max_{\mathbf{v} \in \mathcal{V}} \sup_{\|\mathbf{w}\|_2 \leq B_w} \bar{U}_v(\mathbf{w}) \leq C' B_{xv}^2 \frac{d \log(8g = \nu) + \log(1 - \epsilon)}{N};$$

where we use $\log \nu \leq d \log(4g = \nu)$. Therefore, we plug in the definition of ν and B_{xv} to deduce that, if we set

$$C_1 = \frac{16C' B_{xv}^2}{x} \log(8g = \nu); \quad \nu = \frac{128K_x^2}{x}; \quad B_{xv} = K_x(15 + \log(K_x^2 = x));$$

then as long as $N \geq C_1 d$, it holds $\max_{\mathbf{v} \in \mathcal{V}} E[U_v(\mathbf{w})] - U_v(\mathbf{w}) \leq \frac{x B_{xv}}{16}$ with probability at least $1 - \epsilon \exp(-N/C_1)$. This is the desired result. \square

Lemma G.3. Under Assumption B, for $B_{xv} = K_x(15 + \log(K_x^2 = x))$, it holds

$$\inf_{\mathbf{w} \in \mathcal{B}_2(B_w), \mathbf{v} \in \mathcal{S}^{d-1}} E[|f(\mathbf{x}^\top \mathbf{w}) - 2g(\mathbf{x}^\top \mathbf{v})^2| f(\mathbf{x}^\top \mathbf{v})] \geq B_{xv} g] \geq 3 \quad x=4:$$

Proof. For any fixed $\mathbf{w} \in \mathcal{B}_2(B_w); \mathbf{v} \in \mathcal{S}^{d-1}$,

$$\begin{aligned} & E[|f(\mathbf{x}^\top \mathbf{w}) - 2g(\mathbf{x}^\top \mathbf{v})^2| f(\mathbf{x}^\top \mathbf{v})] \geq B_{xv} g] \\ &= E[|f(\mathbf{x}^\top \mathbf{w}) - 2g(\mathbf{x}^\top \mathbf{v})^2| g] \geq E[|f(\mathbf{x}^\top \mathbf{w}) - 2g(\mathbf{x}^\top \mathbf{v})^2| f(\mathbf{x}^\top \mathbf{v})] > B_{xv} g] \\ & \geq E[(\mathbf{x}^\top \mathbf{v})^2 | f(\mathbf{x}^\top \mathbf{v})] > B_{xv} g]: \end{aligned}$$

Because $\mathbf{x} \in \text{SG}(K_x)$, $\mathbf{x}^\top \mathbf{v} \in \text{SG}(K_x)$, and a simple calculation yields

$$E[(\mathbf{x}^\top \mathbf{v})^2 | f(\mathbf{x}^\top \mathbf{v})] > t K_x g] \geq 2K_x^2 (t^2 + 1) \exp(-t^2):$$

Taking $t = 15 + \log(K_x^2 = x)$ gives $E[(\mathbf{x}^\top \mathbf{v})^2 | f(\mathbf{x}^\top \mathbf{v})] > B_{xv} g \quad x=4$, which completes the proof. \square

G.4 Proof of Theorem G.3 (b)

Notice that

$$r \mathbb{E}_N(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (g(\mathbf{h}\mathbf{w}; \mathbf{x}_i) - y_i) \mathbf{x}_i$$

In the following, we pick a minimal $1/2$ -covering of S^{d-1} (so $|V| \leq 5^d$). Then by Lemma G.2, it holds

$$r \mathbb{E}_N(\mathbf{w}) - \mathbb{E}[r \mathbb{E}_N(\mathbf{w})] \leq 2 \max_{\mathbf{v} \in V} \left| \frac{r \mathbb{E}_N(\mathbf{w}); \mathbf{v}}{|\mathcal{Z}|} - \frac{\mathbb{E}[r \mathbb{E}_N(\mathbf{w}); \mathbf{v}]}{|\mathcal{Z}|} \right|$$

$\mathcal{Z} =: X_{\mathbf{v}}(\mathbf{w})$

Fix a $\mathbf{v} \in S^{d-1}$ and set $|\mathcal{Z}| = |V|$. We proceed to bound $\sup_{\mathbf{w}} |r X_{\mathbf{v}}(\mathbf{w})|$ by applying Proposition B.4 to the random process $f(\mathbf{z}; \mathbf{w}) = g(\mathbf{h}\mathbf{w}; \mathbf{x}_i) - y_i$. We need to verify the preconditions:

- (a) With norm $\|\mathbf{w} - \mathbf{w}'\| = k\mathbf{w} - \mathbf{w}'\|_2$, $\log N(\epsilon; B(r); \|\cdot\|) \leq d \log(2A/r\epsilon)$ with constant $A = 2$;
- (b) For $\mathbf{z} = [\mathbf{x}; y]$, we let $f(\mathbf{z}; \mathbf{w}) := (g(\mathbf{h}\mathbf{w}; \mathbf{x}_i) - y) \mathbf{x}_i$, then $f(\mathbf{z}; \mathbf{w}) \in \text{SE}(CK_x K_y)$ for any \mathbf{w} by our assumption on $(\mathbf{x}; y)$;
- (c) For $\mathbf{w}, \mathbf{w}' \in W$, we have $|g(\mathbf{h}\mathbf{w}; \mathbf{x}_i) - g(\mathbf{h}\mathbf{w}'; \mathbf{x}_i)| \leq L_g \|\mathbf{h}\mathbf{w} - \mathbf{h}\mathbf{w}'\|_2$. Hence, because $\mathbf{x} \in \text{SG}(K_x)$, the random variable $g(\mathbf{h}\mathbf{w}; \mathbf{x}_i) - g(\mathbf{h}\mathbf{w}'; \mathbf{x}_i)$ is sub-Gaussian in $\text{SG}(K_x L_g k\mathbf{w} - \mathbf{w}')_2$. Thus, $f(\mathbf{z}; \mathbf{w}) - f(\mathbf{z}; \mathbf{w}')$ is sub-exponential in $\text{SE}(CK_x^2 L_g k\mathbf{w} - \mathbf{w}')_2$.

Therefore, we can apply Proposition B.4 to obtain that with probability $1 - \delta$, it holds

$$\sup_{\mathbf{w}} |r X_{\mathbf{v}}(\mathbf{w})| \leq C' K_x K_y \left(\frac{d \log(2/\delta) + \log(1/\delta)}{N} + \frac{d \log(2/\delta) + \log(1/\delta)}{N} \right);$$

where we denote $\delta = 1 - L_g K_x^2 B_w = K_y$. Setting $\delta = |V|^{-1}$ and taking the union bound over $\mathbf{v} \in V$, we obtain that with probability at least $1 - \delta$,

$$\max_{\mathbf{v} \in V} \sup_{\|\mathbf{w}\|_2 \leq B_w} |r X_{\mathbf{v}}(\mathbf{w})| \leq C' K_x K_y \left(\frac{d \log(10/\delta) + \log(1/\delta)}{N} + \frac{d \log(10/\delta) + \log(1/\delta)}{N} \right);$$

This is the desired result. \square

G.5 Proof of Theorem G.3 (c)

In the following, we condition on (a) holds, i.e. \mathbb{E}_N is μ -strongly-convex and L -smooth over $B_2(B_w)$ with $\mu = \mu_x/g=8$ and $L = 8L_g K_x^2$. We define

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in B_2(B_w)} \mathbb{E}_N(\mathbf{w});$$

Then by standard convex analysis, we have

$$k\mathbf{w} - \mathbf{w}^*\|_2^2 \leq \frac{D}{\mu} (r \mathbb{E}_N(\mathbf{w}) - r \mathbb{E}_N(\mathbf{w}^*)); \mathbf{w} \in B_2(B_w) \implies r \mathbb{E}_N(\mathbf{w}) - r \mathbb{E}_N(\mathbf{w}^*) \leq \frac{D}{\mu} k\mathbf{w} - \mathbf{w}^*\|_2^2$$

Notice that $r \mathbb{E}_N(\mathbf{w}^*) \leq \mu_{\text{stat}}$, we can conclude that

$$k\mathbf{w}\|_2 \leq k\mathbf{w}^*\|_2 + \frac{\mu_{\text{stat}}}{\mu};$$

Recall that we assume $k\mathbf{w}^*\|_2 \leq B_w/4$, we can then consider $E_S := \mu_{\text{stat}} < B_w/4g$. Once E_S holds, our argument above yields $k\mathbf{w}\|_2 < B_w$, which implies $r \mathbb{E}_N(\mathbf{w}) = 0$. Therefore, $\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \mathbb{E}_N(\mathbf{w})$. Further, by Theorem G.3, we can set

$$C_2 := \max \left(2 \frac{32 K_x K_y}{B_w}, 2 \frac{32 K_x K_y}{B_w} \right);$$

so that as long as $N \geq C_2 d$, the event E_S holds with probability at least $1 - \exp(-N/C_2)$. This is the desired result. \square

G.6 Proof of Theorem G.3 (d) & (e)

We first prove Theorem G.3 (d). Notice that

$$r^2 L_p(\mathbf{w}) = \mathbb{E} g'(\mathbf{h}\mathbf{x}; \mathbf{w}) \mathbf{x} \mathbf{x}^\top \quad \mathbb{E} g''(\mathbf{h}\mathbf{x}; \mathbf{w}) \mathbf{x} \mathbf{x}^\top \quad g'(\mathbf{h}\mathbf{x}; \mathbf{w}) \mathbf{x} \mathbf{x}^\top \quad g''(\mathbf{h}\mathbf{x}; \mathbf{w}) \mathbf{x} \mathbf{x}^\top \quad \mathbf{I}_d; \delta \mathbf{w} \geq B_2(B_w):$$

Therefore, L_p is $(\frac{1}{g''})$ -strongly-convex over $B_2(B_w)$. Therefore, because $\mathbf{w}^* \in B_2(B_w)$ is the global minimum of L_p , it holds that for all $\mathbf{w} \in B_2(B_w)$,

$$L_p(\mathbf{w}) - L_p(\mathbf{w}^*) \leq \frac{1}{2} \frac{1}{g''} \|\mathbf{w} - \mathbf{w}^*\|_2^2.$$

By the definition of $\|\cdot\|_{\text{stat}}$, $\|\mathbf{w} - \mathbf{w}^*\|_{\text{stat}} \leq \frac{1}{g''} \|\mathbf{w} - \mathbf{w}^*\|_2$, and hence the proof of Theorem G.3 (d) is completed.

We next prove Theorem G.3 (e), where we assume that $\mathbb{E}[y|\mathbf{x}] = g(\mathbf{h}\mathbf{x}; \theta)$ (which implies $\mathbf{w}^* = \theta$ directly) and $\mathbf{w}_{\text{GLM}} \in B_2(B_w)$. Notice that

$$r L_p(\mathbf{w}) = \mathbb{E} r \mathbb{E}_N(\mathbf{w}) = \mathbb{E}[(g(\mathbf{h}\mathbf{x}; \mathbf{w}) - y)\mathbf{x}] = \mathbb{E}[(g(\mathbf{h}\mathbf{x}; \mathbf{w}) - g(\mathbf{h}\mathbf{w}; \theta))\mathbf{x}];$$

and hence

$$\begin{aligned} \|\mathbf{w}_{\text{GLM}} - \theta\|_{\text{stat}} &= \mathbb{E}[(g(\mathbf{h}\mathbf{x}; \mathbf{w}_{\text{GLM}}) - g(\mathbf{h}\mathbf{w}; \theta)) (\mathbf{h}\mathbf{x}; \mathbf{w}_{\text{GLM}} - \mathbf{h}\mathbf{w}; \theta)] \\ &= \frac{1}{L_g} \mathbb{E} (g(\mathbf{h}\mathbf{x}; \mathbf{w}_{\text{GLM}}) - g(\mathbf{h}\mathbf{w}; \theta))^2. \end{aligned}$$

On the other hand, by the $(\frac{1}{g''})$ -strong-convexity of L_p over $B_2(B_w)$, it holds that

$$\|\mathbf{w}_{\text{GLM}} - \theta\|_{\text{stat}} \leq \frac{1}{g''} \|\mathbf{w}_{\text{GLM}} - \theta\|_2.$$

Finally, using the definition of \mathbf{w}_{GLM} , we have $\|\mathbf{w}_{\text{GLM}} - \theta\|_{\text{stat}} = 0$, and hence $\|\mathbf{w}_{\text{GLM}} - \theta\|_2 = 0$, which completes the proof of Theorem G.3 (e). \square

H Proofs for Section 3.2

H.1 Proof of Theorem 7

Fix $N \geq 0$, $\epsilon > 0$ and $B_w > 0$, and consider any in-context data D such that the precondition of Theorem 7 holds. Recall that

$$L_{\text{lasso}}(\mathbf{w}) := \frac{1}{2N} \sum_{i=1}^N (\mathbf{h}\mathbf{w}; \mathbf{x}_i - y_i)^2 + N \|\mathbf{w}\|_1$$

denotes the lasso regression loss in (ICLasso), so that $\mathbf{w}_{\text{lasso}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} L_{\text{lasso}}(\mathbf{w})$. We further write

$$\mathbb{E}_N^0(\mathbf{w}) := \frac{1}{2N} \sum_{i=1}^N (\mathbf{h}\mathbf{w}; \mathbf{x}_i - y_i)^2; \quad R(\mathbf{w}) := N \|\mathbf{w}\|_1.$$

Note that $r^2 \mathbb{E}_N^0(\mathbf{w}) = \mathbf{X}^\top \mathbf{X} / N$ and thus \mathbb{E}_N^0 is $\frac{1}{N}$ -smooth over \mathbb{R}^d .

Consider the proximal gradient descent algorithm on the ridge loss

$$\mathbf{w}_{\text{PGD}}^{t+1} = \text{prox}_{\mathcal{R}} \left(\mathbf{w}_{\text{PGD}}^t - \eta \mathbb{E}_N^0(\mathbf{w}_{\text{PGD}}^t) \right)$$

with initialization $\mathbf{w}_{\text{PGD}}^0 := \mathbf{0}_d$, learning rate $\eta := \frac{1}{N}$, and number of steps T to be specified later. Similar to the proof of Theorem 4, we can construct a transformer to approximate \mathbf{w}_{GD}^T . Consider $\psi(s; t) = \frac{1}{2}(s - t)^2$ and $R(\mathbf{w}) = N \|\mathbf{w}\|_1$, then $\psi(s; t)$ is $(0; +1; 2; 4)$ -approximable by sum of relus (cf. Definition D.1), and $\text{prox}_{\mathcal{R}}$ is $(0; +1; 4d; 4 + 2/N)$ -approximable by sum of relus (Proposition D.1). Therefore, we can apply Theorem D.2 with the square loss ψ , regularizer R , learning rate η and accuracy parameter ϵ to obtain that there exists a transformer TF with $(T + 1)$

layers, number of heads $M^{(\ell)} = 2$ for all $\ell \geq [L]$, and hidden dimension $D^\ell = 2d$, such that the final output $\mathbf{h}_{N+1}^{(L)} = [\mathbf{X}_{N+1}; \mathbf{y}_{N+1}; \mathbf{w}_{\text{PGD}}^T]$ with $\mathbf{y}_{N+1} = \mathbf{w}_{\text{PGD}}^T \mathbf{X}_{N+1}$. Further, the weight matrices have norm bounds $\|\mathbf{w}_{\text{PGD}}\| \leq 10R + (8 + 2/N)^{-1}$.

By the standard convergence result for proximal gradient descent (Proposition B.3), we have for all $t \geq 1$ that

$$L_{\text{lasso}}(\mathbf{w}_{\text{PGD}}^t) - L_{\text{lasso}}(\mathbf{w}_{\text{lasso}}) \leq \frac{1}{2t} \|\mathbf{w}_{\text{lasso}}\|_2^2.$$

Plugging in $\|\mathbf{w}_{\text{lasso}}\|_2 \leq B_w$ and $T = L - 1 = \lfloor \frac{1}{\epsilon} \rfloor$ finishes the proof. \square

H.2 Sharper convergence analysis of proximal gradient descent for Lasso

Collection of parameters Throughout the rest of this section, we consider fixed $N \geq 1$, $N = \frac{\log d}{\epsilon}$ for $\epsilon > 0$, ϵ fixed (and to be determined), fixed $0 < \epsilon$, and fixed $B_w^2 > 0$. We write $\epsilon := \frac{\log d}{N}$; $s := (B_w^2)^2$, and $\epsilon_N := -\frac{s \log d}{N}$.

Here we present a sharper convergence analysis on the proximal gradient descent algorithm for L_{lasso} under the following well-conditionedness assumption, which will be useful for proving Theorem 8 in the sequel.

Assumption C (Well-conditioned property for Lasso). *We say the (ICLasso) problem is well-conditioned with sparsity s if the following conditions hold:*

1. The $(\epsilon; s)$ -RSC condition holds:

$$\frac{\|\mathbf{X}\mathbf{w}\|_2^2}{N} \geq \|\mathbf{w}\|_2^2 - \frac{\log d}{N} \|\mathbf{w}\|_1^2; \quad \forall \mathbf{w} \in \mathbb{R}^d. \quad (38)$$

Further, $\max(\mathbf{X}^T \mathbf{X}) = N$.

2. The data $(\mathbf{X}; \mathbf{y})$ is ‘‘approximately generated from a s -sparse linear model’’: There exists a $\mathbf{w}_? \in \mathbb{R}^d$ such that $\|\mathbf{X}\mathbf{w}_? - \mathbf{y}\|_2 \leq B_w^2 \|\mathbf{w}_?\|_1$ and for the residue $\mathbf{r} = \mathbf{y} - \mathbf{X}\mathbf{w}_?$,

$$\|\mathbf{X}^T \mathbf{r}\|_\infty \leq \frac{1}{2} N \epsilon_N.$$

3. It holds that $N \geq 32 - s \log d$ (i.e. $32! N \geq 1$).

Assumption C1 imposes the standard restricted strong convexity (RSC) condition for the feature matrix $\mathbf{X} \in \mathbb{R}^{N \times d}$, and Assumption C2 asserts that the data is approximately generated from a sparse linear model, with a bound on the L_∞ norm of the error vector $\mathbf{X}^T \mathbf{r}$. Assumption C is entirely deterministic in nature, and suffices to imply the following convergence result. In the proof of Theorem 8, we show that Assumption C is satisfied with high probability when data is generated from the standard sparse linear model considered therein.

Theorem H.1 (Sharper convergence guarantee for Lasso). *Under Assumption C, for the PGD iterates $\mathbf{w}_{t \geq 0}^t$ on loss function L_{lasso} with stepsize $\eta = 1$ and starting point $\mathbf{w}^0 = \mathbf{0}$, we have $L_{\text{lasso}}(\mathbf{w}^T) - L_{\text{lasso}}(\mathbf{w}_{\text{lasso}}) \leq \epsilon$ for all*

$$T \leq C \frac{(B_w^2)^2}{\epsilon} + \log C \frac{(B_w^2)^2}{\epsilon} + \frac{1}{\epsilon} N \epsilon_N;$$

where C is a universal constant.

The proof can be found in Appendix H.4. Combining Theorem H.1 with the construction in Theorem 7, we directly obtain the following result as a corollary.

Theorem H.2 (In-context Lasso with transformers with sharper convergence). *For any $N; d; s \geq 1$, $0 < \epsilon$, $\epsilon > 0$, $\epsilon > 0$, there exists a L -layer transformer TF with*

$$L = C \frac{1}{\epsilon} + (\log(C \frac{1}{\epsilon}) + \frac{1}{\epsilon} N \epsilon_N); \quad \max_{\ell \in [L]} M^{(\ell)} \leq 2; \quad \max_{\ell \in [L]} D^{(\ell)} \leq 2d;$$

Proof. We follow the notation in the proof of Lemma H.1. By (39), we have

$$0 \leq \frac{1}{2N} \mathbf{X}^T \mathbf{X} \mathbf{k}_2^2 - \frac{N}{2} (3s - s k_1 - k - s c k_1) + \mathfrak{L}_{\text{lasso}}(\mathbf{w}) - \mathfrak{L}_{\text{lasso}}(\mathbf{w}_?);$$

and hence $k - k_1 \leq 4 \frac{\rho_-}{s} s k - k_2 + \frac{2\text{gap}}{N}$ due to $\mathfrak{L}_{\text{lasso}}(\mathbf{w}) - \mathfrak{L}_{\text{lasso}}(\mathbf{w}_?) \leq \text{gap}$. On the other hand, by the RSC condition (38), it holds that

$$\frac{k \mathbf{X}^T \mathbf{X} \mathbf{k}_2^2}{N} \leq k - k_2^2 \leq \frac{\log d}{N} k - k_1^2;$$

Therefore, we have

$$\begin{aligned} k - k_2^2 &\leq 3 \frac{\rho_-}{N s} s k - k_2 + \frac{\log d}{N} k - k_1^2 + 2\text{gap} \\ &\leq 3 \frac{\rho_-}{N s} s k - k_2 + \frac{\log d}{N} 4 \frac{\rho_-}{s} s k - k_2 + \frac{2\text{gap}}{N}^2 + 2\text{gap} \\ &\leq \frac{5s}{N} \frac{k^2}{2} + \frac{1}{6} k - k_2^2 + \frac{20s \log d}{2N} k - k_2^2 + \frac{20 \log d}{N} \text{gap}^2 + 2\text{gap}; \end{aligned}$$

where the last inequality uses AM-GM inequality and Cauchy inequality. Notice that $\frac{20s \log d}{N} \leq \frac{2}{3}$, we now derive that

$$k - k_2^2 \leq \frac{30s}{2} \frac{k^2}{N} + \frac{120 \log d}{2N} \text{gap}^2 + 12\text{gap}.$$

Plugging in $N = \frac{c \log d}{\frac{\log d}{N}}$ completes the proof. The corollary follows immediately by letting $\mathbf{w} = \mathbf{w}_{\text{lasso}}$ in above proof (hence $\text{gap} = 0$). \square

Lemma H.2 (Growth). *It holds that*

$$\frac{1}{2N} \mathbf{X}^T \mathbf{X} (\mathbf{w} - \mathbf{w}_{\text{lasso}})^2 \leq \mathfrak{L}_{\text{lasso}}(\mathbf{w}) - \mathfrak{L}_{\text{lasso}}(\mathbf{w}_{\text{lasso}}); \quad \forall \mathbf{w};$$

Proof. For simplicity we denote $\mathbf{w}_{\text{lasso}} := \mathbf{w}_{\text{lasso}}$. By the first order optimality condition, it holds that

$$0 \geq \frac{1}{N} \mathbf{X}^T (\mathbf{X} \mathbf{w}_{\text{lasso}} - \mathbf{y}) + \partial R(\mathbf{w}_{\text{lasso}});$$

where we write $R(\mathbf{w}) := \frac{1}{N} \mathbf{X}^T \mathbf{X} \mathbf{w} - k_1$. Then by the convexity of R , we have

$$\begin{aligned} R(\mathbf{w}) - R(\mathbf{w}_{\text{lasso}}) &\geq \langle \partial R(\mathbf{w}_{\text{lasso}}); \mathbf{w} - \mathbf{w}_{\text{lasso}} \rangle = \frac{1}{N} \mathbf{X}^T (\mathbf{X} \mathbf{w}_{\text{lasso}} - \mathbf{y}); \mathbf{w} - \mathbf{w}_{\text{lasso}} \\ &= \frac{1}{N} \langle \mathbf{X} \mathbf{w}_{\text{lasso}} - \mathbf{y}; \mathbf{X} \mathbf{w} - \mathbf{y} \rangle - \langle \mathbf{X} \mathbf{w}_{\text{lasso}} - \mathbf{y}; \mathbf{w} - \mathbf{w}_{\text{lasso}} \rangle \\ &= \frac{1}{2N} \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{y} k_2^2 + \frac{1}{2N} \mathbf{X}^T \mathbf{X} \mathbf{w}_{\text{lasso}} - \mathbf{y} k_2^2 + \frac{1}{2N} \mathbf{X}^T \mathbf{X} (\mathbf{w} - \mathbf{w}_{\text{lasso}})^2 k_2^2; \end{aligned}$$

Rearranging completes the proof. \square

H.4 Proof of Theorem H.1

For the simplicity of presentation, we write $\mathbf{w}_{\text{lasso}} = \mathbf{w}_{\text{lasso}}$ and we denote $\text{gap}^t := \mathfrak{L}_{\text{lasso}}(\mathbf{w}^t) - \mathfrak{L}_{\text{lasso}}(\mathbf{w}_{\text{lasso}})$.

By Lemma H.1, we have $k \mathbf{w}^t - \mathbf{w}_? k_1 \leq 4 \frac{\rho_-}{s} k \mathbf{w}^t - \mathbf{w}_? k_2 + \frac{2\text{gap}^t}{N}$, which implies

$$\mathbf{w}^t - \mathbf{w}_{\text{lasso}} \leq \mathbf{w}_? \frac{1}{s} + k \mathbf{w}_{\text{lasso}} - \mathbf{w}_? k_1 \leq 4 \frac{\rho_-}{s} \mathbf{w}^t - \mathbf{w}_{\text{lasso}} \frac{1}{s} + 8 \frac{\rho_-}{s} k \mathbf{w}_{\text{lasso}} - \mathbf{w}_? k_2 + \frac{2\text{gap}^t}{N};$$

We denote $N = \frac{2 \log d}{N}$. Using the assumption that \mathbf{X} is $(\cdot; \cdot)$ -RSC, we obtain that

$$\frac{1}{N} \mathbf{X}(\mathbf{w}^t - \mathbf{w}_{\text{lasso}})^2 \leq \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 + 640s \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 + \frac{40}{2} (\text{gap}^t)^2 + \frac{40}{2} (\text{gap}^t)^2 ;$$

Thus, as long as $N \geq \frac{30 \cdot 2s \log d}{\cdot}$, we have

$$\frac{1}{3} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 \leq \frac{1}{N} \mathbf{X}(\mathbf{w}^t - \mathbf{w}_{\text{lasso}})^2 + 640s \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 + \frac{40}{2} \frac{N}{N} (\text{gap}^t)^2 + 2\text{gap}^t + 40 \frac{1}{N} (\text{gap}^t)^2 + 640s \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 ;$$

where the last inequality follows from Lemma H.2 and the definition of N .

We define $\eta_{\text{stat}} := 640s \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2$, $T_0 := 10 \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2$. By Proposition B.3(3), it holds that for $t \geq T_0$,

$$\text{gap}^t \leq \frac{1}{2t} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 \leq \frac{1}{2T_0} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 = \frac{1}{20} ;$$

Then for all $t \geq T_0 + 1$, we have (the second line below uses Proposition B.3(2))

$$\frac{1}{3} \mathbf{w}^{t+1} \mathbf{w}_{\text{lasso}}^2 \leq 4\text{gap}^{t+1} + \eta_{\text{stat}} \frac{1}{2} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 + \frac{1}{3} \mathbf{w}^{t+1} \mathbf{w}_{\text{lasso}}^2 + \eta_{\text{stat}} ;$$

$$\Rightarrow \mathbf{w}^{t+1} \mathbf{w}_{\text{lasso}}^2 \leq \frac{3\eta_{\text{stat}}}{1 + \frac{1}{6}} \frac{1}{2} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 + \frac{3\eta_{\text{stat}}}{1} ;$$

Therefore, for $t \geq T_0 + 1$,

$$\mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 \leq \exp\left(-\frac{1}{12}(t - T_0)e + 1\right) \mathbf{w}^{\lceil T_0 \rceil} \mathbf{w}_{\text{lasso}}^2 + \frac{3\eta_{\text{stat}}}{2} ;$$

$$\leq \exp\left(-\frac{1}{8}(t - T_0) \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 + \frac{3\eta_{\text{stat}}}{2}\right) ;$$

where the last inequality follows from Proposition B.3(2). Further, by Proposition B.3(3), we have

$$\text{gap}^{t+k} \leq \frac{1}{2k} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 \leq \frac{1}{2k} \exp\left(-\frac{1}{8}(t - T_0) \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 + \frac{3\eta_{\text{stat}}}{2}\right) ; \quad \forall t \geq T_0 + 1, k \geq 0 ;$$

Hence, we can conclude that $\text{gap}^T \leq \frac{1}{2k}$ for all T such that

$$T \geq 10 \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 + 8 \log \frac{\frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2}{\frac{1}{2k}} + \frac{3\eta_{\text{stat}}}{\frac{1}{2k}} + 1 ;$$

Now, by Proposition H.1, it holds that $\frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 \leq 10 \frac{1}{2} \frac{s \log d}{N}$, and hence

$$\frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 \leq 2 \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 + 2 \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 \leq 2(B_w^2)^2 + \frac{20}{2} \frac{s \log d}{N} ;$$

Plugging in our definition of

$$N = \frac{\log d}{N} ; \quad \eta_{\text{stat}} := 400s \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 ; \quad \frac{1}{N} \mathbf{w}^t \mathbf{w}_{\text{lasso}}^2 = \frac{s \log d}{N} + 1$$

completes the proof. \square

H.5 Proof of Theorem 8

In this section, we present the proof of Theorem 8 based on Theorem H.2. We begin by recalling the following RSC property of a Gaussian random matrix [87, Theorem 7.16], a classical result in the high-dimensional statistics literature.

Proposition H.2 (RSC for Gaussian random design). *Suppose that $\mathbf{X} = [\mathbf{x}_1; \dots; \mathbf{x}_N]^\top \in \mathbb{R}^{N \times d}$ is a random matrix with each row \mathbf{x}_i being i.i.d. samples from $\mathcal{N}(0; \Sigma)$. Then there are universal constants $c_1 = \frac{1}{8}; c_2 = 50$ such that with probability at least $1 - \frac{e^{-N/32}}{1 - e^{-N/32}}$,*

$$\frac{k\mathbf{X}\mathbf{w}k_2^2}{N} \leq c_1 k\mathbf{w}k^2 + c_2 \left(\frac{\log d}{N}\right) k\mathbf{w}k_1^2; \quad \forall \mathbf{w} \in \mathbb{R}^d; \quad (40)$$

where $\lambda_{\max}(\Sigma) = \max_{i \in [d]} \Sigma_{ii}$ is the maximum of diagonal entries of Σ .

Fix a parameter β (which we will specify in proof) and a large universal constant C_0 . Let us set

$$\beta = c_1 = \frac{1}{8}; \quad \beta = 8(1 + (d=N)); \quad \beta = c_2 = \frac{1}{8};$$

$$B_x = C_0 \frac{\beta}{d \log(N-1)}; \quad B_y = C_0 (\beta^2 + \frac{\beta}{d \log(N-1)});$$

Similar to the proof of Corollary 6 (Appendix F.4), we consider the following good events (where $\mathbf{y} = \mathbf{X}\mathbf{w} + \mathbf{y}$)

$$E_w = \mathbb{P} \left[\max_{\|\mathbf{w}\|_1 \leq N} (\mathbf{X}^\top \mathbf{X} - N) \mathbf{w} \leq \beta \text{ and } \mathbf{X} \text{ is } \left(\frac{\beta}{N}; \beta\right)\text{-RSC} \right];$$

$$E_r = \mathbb{P} \left[\|\mathbf{X}^\top \mathbf{y}\|_\infty \leq 4 \frac{\beta}{N \log(4d)} \right];$$

$$E_b = \mathbb{P} \left[\|\mathbf{x}_i\|_2 \leq B_x; \|\mathbf{y}_i\|_2 \leq B_y; \forall i \in [N] \right];$$

$$E_{b;N+1} = \mathbb{P} \left[\|\mathbf{x}_{N+1}\|_2 \leq B_x; \|\mathbf{y}_{N+1}\|_2 \leq B_y \right];$$

and we define $E := E_w \setminus E_r \setminus E_b \setminus E_{b;N+1}$.

Furthermore, we choose $\beta > 0$ that correspond to the choice $\beta = 8 \frac{\log(4d)}{N}$, and we also assume $N \geq \frac{32c_2}{c_1} s \log d$. Then, Assumption C holds on the event E .

Therefore, we can apply Theorem H.2 with $\beta = \beta/N$, which implies that there exists a L -layer transformer \mathcal{T} such that its prediction $\hat{\mathbf{y}}_{N+1} := \text{read}_y(\mathcal{T}^0(\mathbf{H}))$, so that under the good event E we have $\hat{\mathbf{y}}_{N+1} = \text{clip}_{B_y}(\mathbf{h}_{N+1}; \mathbf{w})$, where

$$L_{\text{lasso}}(\mathbf{w}) = L_{\text{lasso}}(\mathbf{w}_{\text{lasso}}) + \beta/N;$$

In the following, we show that \mathcal{T} is indeed the desired transformer (similarly to the proof in Appendix F.4). Consider the conditional prediction error

$$\mathbb{E} \left[(\hat{\mathbf{y}}_{N+1} - \mathbf{y}_{N+1})^2 \mid D \right] = \mathbb{E} \left[1FEg(\hat{\mathbf{y}}_{N+1} - \mathbf{y}_{N+1})^2 \mid D \right] + \mathbb{E} \left[1FE^c g(\hat{\mathbf{y}}_{N+1} - \mathbf{y}_{N+1})^2 \mid D \right];$$

and we analyze these two parts separately under the good event $E_0 := E_w \setminus E_r \setminus E_b$ of D .

Part I. We first note that

$$\mathbb{E} \left[1FEg(\hat{\mathbf{y}}_{N+1} - \mathbf{y}_{N+1})^2 \mid D \right] \leq \mathbb{E} \left[1FEg(\text{clip}_{B_y}(\mathbf{h}_{N+1}; \mathbf{w}) - \mathbf{y}_{N+1})^2 \mid D \right]$$

$$\leq \mathbb{E} \left[1FEg(\mathbf{h}_{N+1}; \mathbf{w}) - \mathbf{y}_{N+1})^2 \mid D \right];$$

where the inequality is because $\mathbf{y}_{N+1} \in [B_y; B_y]$ under the good event E . Notice that by our construction, under the good event E , $\mathbf{w} = \mathbf{w}(D)$ depends only on the dataset D (because it is the $(L-1)$ -th iterate of PGD on (ICLasso) problem). Applying Proposition H.1 to $\mathbf{w}(D)$ and using the definition of β/N and our choice of β/N , we obtain that (under E_0)

$$k\mathbf{w}(D)k_2^2 \leq C \left(\frac{s}{N} + \beta^2/N + \beta/N \right) = O \left(\frac{s \log(d)}{N} \right);$$

Therefore, under E_0 ,

$$\mathbb{E} \left[1FEg(\mathbf{h}_{N+1}; \mathbf{w}(D) - \mathbf{y}_{N+1})^2 \mid D \right] = \mathbb{E} \left[1FEg(\mathbf{h}_{N+1}; \mathbf{w}(D) - \mathbf{y}_{N+1})^2 \mid D \right]$$

$$= \mathbb{E} \left[(\mathbf{h}_{N+1}; \mathbf{w}(D) - \mathbf{y}_{N+1})^2 \mid D \right]$$

$$= \mathbb{E} \left[(\mathbf{h}_{N+1}; \mathbf{w}(D) - \mathbf{h}_{N+1}; \mathbf{w}(D)) + (\mathbf{h}_{N+1}; \mathbf{w}(D) - \mathbf{h}_{N+1}; \mathbf{w}(D) - \mathbf{y}_{N+1})^2 \mid D \right] + \beta^2$$

$$= k\mathbf{w}(D)k_2^2 + \beta^2$$

$$= \beta^2 + O \left(\frac{s \log(d)}{N} \right);$$

Part II. Notice that under good event E_0 , the bad event E^c holds if and only if $E_{b;N+1}^c$ holds, and hence

$$E \mathbb{1}_{E^c} g(y_{N+1}, y_{N+1})^2 D = E \mathbb{1}_{E_{b;N+1}^c} g(y_{N+1}, y_{N+1})^2 D \leq \frac{E[(y_{N+1}, y_{N+1})^4]}{P(E_{b;N+1}^c)}.$$

With a large enough constant C_0 , we clearly have $P(E_{b;N+1}^c) \geq (1-N)^{10}$. Further, a simple calculation yields

$$E[(y_{N+1}, y_{N+1})^4] \leq 8E[y_{N+1}^4 + y_{N+1}^4] \leq 8B_y^4 + 8E[y_{N+1}^4] \leq 16B_y^4,$$

where the last inequality is because the marginal distribution of y_{N+1} is simply $N(0; \sigma^2 + kw_{\gamma}k_2^2)$. Combining these yields

$$E \mathbb{1}_{E^c} g(y_{N+1}, y_{N+1})^2 D \leq O\left(\frac{B_y^4}{N^5}\right) \leq O\left(\frac{((B_w^2)^2 + \sigma^2) \log(1-N)}{N^4}\right).$$

Therefore, choosing $\delta = \min\{f; \frac{1}{B_w^2}g\}$ is enough for our purpose, and under such choice of δ ,

$$E \mathbb{1}_{E^c} g(y_{N+1}, y_{N+1})^2 D \leq O\left(\frac{\delta^2}{N^4}\right).$$

Conclusion. Combining the inequalities above, we can conclude that under E_0 ,

$$E[(y_{N+1}, y_{N+1})^2] D \leq \delta^2 + O\left(\frac{\log(1-N)}{N}\right).$$

It remains to show that $P(E_0) \geq 1 - \epsilon$. By Proposition H.2, Lemma B.2 and Lemma B.4, we have

$$P(E_w) \leq 3 \exp(-N/32); \quad P(E_r) \leq \frac{\epsilon}{2}; \quad P(E_b) \leq \frac{\epsilon}{4}.$$

Therefore, as long as $N \geq 32 \log(1/\epsilon)$, we have $P(E_0) \geq 1 - \epsilon$. This completes the proof. \square

We also remark that in the construction above,

$$R = O\left((B_w^2 + \sigma^2)^{p-1} d \log(N(1 + B_w^2))\right);$$

which would be useful for bounding $\mathbb{J} \mathbb{J}$.

I Proofs for Section 4

I.1 Proof of Proposition 10

We begin by restating Proposition 10 into the following version, which contains additional size bounds on \mathcal{D}_{val} .

Theorem I.1 (Full statement of Proposition 10). *Suppose that for*

$$\mathbb{D}_{\text{val}}(f) := \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{val}}} \sigma(f(\mathbf{x}_i); y_i);$$

$\sigma(\cdot; \cdot)$ is $(\epsilon; R; M; C)$ -approximable by sum of relus (Definition D.1). Then there exists a 3-layer transformer TF with

$$\max_{\epsilon \in [3]} M^{(\epsilon)} \leq (M+3)K; \quad \max_{\epsilon \in [3]} D^{(\epsilon)} \leq K^2 + K + 1; \quad \mathbb{J} \mathbb{J} \leq \frac{2NKC}{|\mathcal{D}_{\text{val}}|} + 3\epsilon^{-1} + 7KR;$$

that maps

$$\mathbf{h}_i = [\epsilon; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbf{0}_{K+1}; 1; t_i] \quad \mathbf{h}'_i = [\epsilon; \mathbb{P}(\mathbf{x}_i); 1; t_i]; \quad i \in [N+1];$$

where the predictor $\mathbb{P}: \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex combination of f_k : $\mathbb{D}_{\text{val}}(f_k) \leq \min_{k \in [K]} \mathbb{D}_{\text{val}}(f_k) + g$. As a corollary, for any convex risk $L: (\mathbb{R}^d \rightarrow \mathbb{R}) \rightarrow \mathbb{R}$, \mathbb{P} satisfies

$$L(\mathbb{P}) \leq \min_{k \in [K]} L(f_k) + \max_{k \in [K]} \mathbb{D}_{\text{val}}(f_k) \leq L(f_k) + g.$$

To prove Theorem I.1, we first state and prove the following two propositions.

Proposition I.1 (Evaluation layer). *There exists a 1-layer transformer TF with MK heads and $\prod_{j=1}^M 3R + 2NK C = jD_{\text{val}}/j$ such that for all \mathbf{H} such that $\max_i f_j y_j/g \leq R; \max_{i,k} f_k f_k(\mathbf{x}_i)/g \leq R$, TF maps*

$$\begin{aligned} \mathbf{h}_i &= [\mathbf{x}_i; y_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbf{0}_{K+1}; 1; t_i] \\ \mathbf{h}'_i &= [\mathbf{x}_i; y_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbb{E}_{\text{val}}(f_1); \dots; \mathbb{E}_{\text{val}}(f_K); 0; 1; t_i]; \quad i \geq [N + 1]; \end{aligned}$$

where $\mathbb{E}_{\text{val}}(\cdot)$ is a functional such that $\max_k \mathbb{E}_{\text{val}}(f_k) = \mathbb{E}_{\text{val}}(f_k)$.

Proof of Proposition I.1. As \cdot is $(\cdot; R; M; C)$ -approximable by sum of relus, there exists a function $g: \mathbb{R}^2 \rightarrow \mathbb{R}$ of form

$$g(s; t) = \sum_{m=1}^M c_m (a_m s + b_m t + d_m) \quad \text{with} \quad \sum_{m=1}^M j c_m j \leq C; \quad |a_m| + |b_m| + |d_m| \leq 1; \quad 8m \leq [M];$$

such that $\sup_{(s;t) \in [-R; R]^2} |g(s; t) - \cdot(s; t)| \leq \epsilon$. We define

$$\mathbb{E}_{\text{val}}(f) := \frac{1}{jD_{\text{val}}/j} \sum_{(\mathbf{x}_i; y_i) \in \mathcal{D}_{\text{val}}} g(f(\mathbf{x}_i); y_i);$$

Next, for every $m \leq [M]$ and $k \leq [K]$, we define matrices $\mathbf{Q}_{m;k}; \mathbf{K}_{m;k}; \mathbf{V}_{m;k} \in \mathbb{R}^{D \times D}$ such that for all $i; j \leq [N + 1]$,

$$\mathbf{Q}_{m;k} \mathbf{h}_i = \begin{bmatrix} a_m \\ b_m \\ d_m \\ 0 \end{bmatrix}; \quad \mathbf{K}_{m;k} \mathbf{h}_j = \begin{bmatrix} f_k(\mathbf{x}_j) \\ y_j \\ 1 \\ R(1 + t_j) \\ 0 \end{bmatrix}; \quad \mathbf{V}_{m;k} \mathbf{h}_j = \frac{(N + 1)c_m}{jD_{\text{val}}/j} \mathbf{e}_{D - (K - k) - 3}$$

where $\mathbf{e}_s \in \mathbb{R}^D$ is the vector with s -th entry being 1 and others being 0. As the input has structure $\mathbf{h}_i = [\mathbf{x}_i; y_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbf{0}_{K+1}; 1; t_i]$, these matrices indeed exist, and further it is straightforward to check that they have norm bounds

$$\max_{m \in [M]; k \in [K]} \|\mathbf{Q}_{m;k}\|_{\text{op}} \leq 3; \quad \max_{m \in [M]; k \in [K]} \|\mathbf{K}_{m;k}\|_{\text{op}} \leq 2 + R; \quad \max_{m \in [M]; k \in [K]} \|\mathbf{V}_{m;k}\|_{\text{op}} \leq \frac{K(N + 1)C}{jD_{\text{val}}/j}.$$

Now, for every $i; j \leq [N + 1]$, we have

$$\begin{aligned} (\mathbf{h} \mathbf{Q}_{m;k} \mathbf{h}_i; \mathbf{K}_{m;k} \mathbf{h}_j) &= (a_m f_k(\mathbf{x}_i) + b_m y_j + d_m - 2R(1 + t_j)) \\ &= a_m \mathbf{w}^\top \mathbf{x}_i + b_m y_j + d_m - 1ft_j = 1g; \end{aligned}$$

where the last equality follows from the bound $|a_m f_k(\mathbf{x}_i) + b_m y_j + d_m - R(ja_m + jb_m)| \leq d_m \leq 2R$, so that the above relu equals 0 if $t_j = 0$. Therefore, for each $i \leq [N + 1]$ and $k \leq [K]$,

$$\begin{aligned} & \sum_{m=1}^M (\mathbf{h} \mathbf{Q}_{m;k} \mathbf{h}_i; \mathbf{K}_{m;k} \mathbf{h}_j) \mathbf{V}_{m;k} \mathbf{h}_j \\ &= \sum_{m=1}^M c_m (a_m \mathbf{w}^\top \mathbf{x}_i + b_m y_j + d_m - 1ft_j) \frac{(N + 1)}{jD_{\text{val}}/j} = 1g \mathbf{e}_{D - (K - k) - 3} \\ &= g(f_k(\mathbf{x}_i); y_j) \frac{(N + 1)}{jD_{\text{val}}/j} 1ft_j = 1g \mathbf{e}_{D - (K - k) - 3}; \end{aligned}$$

Thus letting the attention layer $\text{Attn}_{(m;k)} = f(\mathbf{V}_{m;k}; \mathbf{Q}_{m;k}; \mathbf{K}_{m;k})g_{(m;k) \in [M] \times [K]}$, we have

$$\hat{\mathbf{h}}_i = [\text{Attn}(\mathbf{H})]_i = \mathbf{h}_i + \frac{1}{N + 1} \sum_{j=1}^{N+1} \sum_{m;k} (\mathbf{h} \mathbf{Q}_{m;k} \mathbf{h}_i; \mathbf{K}_{m;k} \mathbf{h}_j) \mathbf{V}_{m;k} \mathbf{h}_j$$

$$\begin{aligned}
&= \mathbf{h}_i + \frac{1}{jD_{\text{val}}^j} \prod_{j=1}^{K+1} \prod_{k=1}^K g(f_k(\mathbf{x}_j); y_j) \mathbf{1}_{f_j} = \mathbf{1}_{g_{D-(K-k)-3}} \\
&= \mathbf{h}_i + \prod_{k=1}^K \frac{1}{jD_{\text{val}}^j} \prod_{(\mathbf{x}_j; y_j) \in \mathcal{D}_{\text{val}}} g(f_k(\mathbf{x}_j); y_j) \mathbf{e}_{D-(K-k)-3} \\
&= \mathbf{h}_i + \prod_{k=1}^K \mathbb{E}_{\text{val}}(f_k) \mathbf{e}_{D-(K-k)-3} \\
&= [\mathbf{x}_i; y_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbf{0}_{K+1}; 1; t_i] + [\mathbf{0}_{D-K-3}; \mathbb{E}_{\text{val}}(f_1); \dots; \mathbb{E}_{\text{val}}(f_K); \mathbf{0}; \mathbf{0}; \mathbf{0}] \\
&= [\mathbf{x}_i; y_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbb{E}_{\text{val}}(f_1); \dots; \mathbb{E}_{\text{val}}(f_K); \mathbf{0}; 1; t_i]; \quad i \geq [N+1];
\end{aligned}$$

This is the desired result. \square

Proposition I.2 (Selection layer). *There exists a 3-layer transformer TF with*

$$\max_{\ell \in [3]} M^{(\ell)} \leq 2K + 2; \quad \max_{\ell \in [3]} D^{(\ell)} \leq K^2 + K + 1; \quad \prod_{\ell \in [3]} \prod_{\ell \in [3]} \ell^{-1} + 3KR + 2;$$

such that TF maps

$$\begin{aligned}
\mathbf{h}_i &= [\mathbf{x}_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); L_1; \dots; L_K; \mathbf{0}; 1; t_i] \\
\mathbf{h}'_i &= [\mathbf{x}_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \dots; \dots; \mathbb{P}(\mathbf{x}_i); 1; t_i]; \quad i \geq [N+1];
\end{aligned}$$

where $\mathbb{P} = \prod_{k=1}^K \alpha_k f_k$ is an aggregated predictor, where the weights $\alpha_1; \dots; \alpha_K \geq 0$ are functions only on $L_1; \dots; L_K$ such that

$$\prod_{k=1}^K \alpha_k = 1; \quad \alpha_k > 0 \text{ only if } L_k = \min_{k' \in [K]} L_{k'} + \dots;$$

Proof of Proposition I.2. We construct a \mathbf{h} which is a composition of 2 MLP layers followed by an attention layer ((1), (2), (3)).

Step 1: construction of $\mathbf{W}_1^{(1)}$. We consider matrix $\mathbf{W}_1^{(1)}$ that maps

$$\begin{aligned}
\mathbf{h} &= [D-K-3; L_1; \dots; L_K; \dots; \dots] \\
\mathbf{W}_1^{(1)} \mathbf{h} &= [L_1 \ L_2; \dots; L_1 \ L_K; \dots; L_K \ L_{K-1}; L_1; L_1; \dots; L_K; L_K];
\end{aligned}$$

i.e. $\mathbf{W}_1^{(1)} \mathbf{h}$ is a $K^2 + K$ dimensional vector so that its entry contains $f_{L_k} \prod_{l \in [K]} g_{k;l \in [K]}$ and $f_{L_k}; \prod_{k \in [K]} g_{k \in [K]}$. Clearly, such $\mathbf{W}_1^{(1)}$ exists and can be chosen so that $\mathbf{W}_1^{(1)} \text{ op} \leq 2K$. We then

consider a matrix $\mathbf{W}_2^{(1)}$ that maps

$$(\mathbf{W}_1^{(1)} \mathbf{h}) \mathbf{W}_2^{(1)} (\mathbf{W}_1^{(1)} \mathbf{h}) = [\mathbf{0}_{D-K-3}; c_1 \ L_1; \dots; c_K \ L_K; \mathbf{0}_3] \in \mathbb{R}^D;$$

where $c_k = c_k(L) := \prod_{l \neq k} (L_k - L_l)$. Notice that

$$c_k - L_k = (L_k) + \prod_{l \neq k} (L_k - L_l);$$

and hence such $\mathbf{W}_2^{(1)}$ exists and can be chosen so that $\mathbf{W}_2^{(1)} \text{ op} \leq K + 1$. We set $\mathbf{W}_{\text{mlp}}^{(1)} = (\mathbf{W}_1^{(1)}; \mathbf{W}_2^{(1)})$, then MLP $\mathbf{W}_{\text{mlp}}^{(1)}$ maps \mathbf{h}_i to

$$\mathbf{h}_i^{(1)} = [\mathbf{x}_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); c_1; \dots; c_K; \mathbf{0}; 1; t_i];$$

The basic property of $f_{c_k} g_{k \in [K]}$ is that, if $c_k = \dots$, then $L_k = \min_{k' \in [K]} L_{k'} + \dots$.

Step 2: construction of $\mathbf{W}_1^{(2)}$. We consider matrix $\mathbf{W}_1^{(2)}$ that maps

$$\mathbf{h} = [D-K-3; c_1; \dots; c_K; \dots; 1; \dots]$$

$$\mathbb{W}_1^{(2)} \mathbf{h} = [1 \quad -1c_1; c_1; c_1; \dots; 1 \quad -1c_K; c_K; c_K] \in \mathbb{R}^{3K};$$

and $\mathbf{W}_1^{(2)}$ can be chosen so that $\mathbf{W}_1^{(2)} \in \mathbb{R}^{(K+1) \times (K+1)}$. We then consider a matrix $\mathbf{W}_2^{(2)}$ that maps

$$(\mathbf{W}_1^{(2)} \mathbf{h}) \in \mathbb{R}^{(K+1)} \quad \mathbf{W}_2^{(2)} \quad (\mathbf{W}_1^{(1)} \mathbf{h}) = [\mathbf{0}_{D-K-3}; (1 \quad -1c_1) \quad c_1; \dots; (1 \quad -1c_K) \quad c_K; \mathbf{0}_3] \in \mathbb{R}^D;$$

which exists and can be chosen so that $\mathbf{W}_2^{(2)} \in \mathbb{R}^{(D-K-3) \times (K+1)}$. We set $\mathbb{M}_{\text{mlp}}^{(2)} = (\mathbf{W}_1^{(2)}; \mathbf{W}_2^{(2)})$, then MLP $\mathbb{M}_{\text{mlp}}^{(2)}$ maps $\mathbf{h}_i^{(1)}$ to

$$\mathbf{h}_i^{(2)} = [f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); u_1; \dots; u_K; 0; 1; t_i];$$

where $u_k = (1 - 1c_k) \delta_k \in [0, 1]$. Clearly, $u_k \in [0, 1]$, and $u_k > 0$ if and only if $c_k < 1$.

Step 3: construction of $\mathbb{A}_{\text{attn}}^{(3)}$. We define

$$A_{k,0} = (1 \quad u_1); \quad A_{k,1} = (1 \quad u_1 \quad \dots \quad u_{k-1}) \quad (1 \quad u_1 \quad \dots \quad u_k) \delta_k \quad 2:$$

Clearly, $A_{k,0} \geq 0$, and $A_{k,1} \geq 0$. Further,

$$A_{k,1} > 0 \iff u_k > 0 \iff c_k < 1 \iff L_k = \min_{k' \in [K]} L_{k'} + 1;$$

Therefore, it remains to construct $\mathbb{A}_{\text{attn}}^{(3)}$ that implements $\mathbb{A} = \sum_{k=1}^K A_{k,1} f_k$ based on $[\mathbf{h}_i^{(2)}]_i$. Notice that

$$\mathbb{A}(\mathbf{x}_i) = (1 \quad f_1(\mathbf{x}_i) + \sum_{k=1}^{K-1} (1 \quad u_1 \quad \dots \quad u_{k-1}) (f_k(\mathbf{x}_i) \quad f_{k-1}(\mathbf{x}_i))) \quad (41)$$

$$(1 \quad u_1 \quad \dots \quad u_K) f_K(\mathbf{x}_i);$$

and hence we construct $\mathbb{A}_{\text{attn}}^{(3)}$ as follows: for every $k \in [K+1]$ and $w \in \{0, 1\}$, we define matrices $\mathbf{Q}_{k,w}; \mathbf{K}_{k,w}; \mathbf{V}_{k,w} \in \mathbb{R}^{D \times D}$ such that for all $k \in [K+1]$

$$\mathbf{Q}_{k,0} \mathbf{h}_i^{(2)} = \begin{pmatrix} f_k(\mathbf{x}_i) + R \\ \mathbf{0} \end{pmatrix} \mathbf{1}_k; \quad \mathbf{Q}_{k,1} \mathbf{h}_i^{(2)} = \begin{pmatrix} f_{k-1}(\mathbf{x}_i) + R \\ \mathbf{0} \end{pmatrix} \mathbf{1}_k;$$

$$\mathbf{K}_{k,0} \mathbf{h}_j^{(2)} = \mathbf{K}_{k,1} \mathbf{h}_j^{(2)} = \begin{pmatrix} 0 & u_1 & \dots & u_{k-1} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & u_{k-1} & \dots & 0 \end{pmatrix}; \quad \mathbf{V}_{k,0} \mathbf{h}_j^{(2)} = \mathbf{e}_{D-2} = \mathbf{V}_{k,1} \mathbf{h}_j^{(2)};$$

for all $i, j \in [N+1]$, where we understand $f_0 = f_{K+1} = 0$ and $\mathbf{1}_k$ is the k -dimensional vector with all entries being 1. By the structure of $\mathbf{h}_i^{(2)}$, these matrices indeed exist, and further it is straightforward to check that they have norm bounds

$$\max_{k \in [K+1]; w \in \{0,1\}} \|\mathbf{Q}_{k,w}\|_{\text{op}} \leq KR; \quad \max_{k \in [K+1]; w \in \{0,1\}} \|\mathbf{K}_{k,w}\|_{\text{op}} \leq 1; \quad \max_{k \in [K+1]; w \in \{0,1\}} \|\mathbf{V}_{k,w}\|_{\text{op}} \leq 2K + 2;$$

Now, for every $i, j \in [N+1], k \in [K+1]; w \in \{0,1\}$, we have

$$\mathbf{Q}_{k,w} \mathbf{h}_i^{(2)}; \mathbf{K}_{k,w} \mathbf{h}_j^{(2)} = ((1 \quad u_1 \quad \dots \quad u_{k-1}) (f_{k-w}(\mathbf{x}_i) + R))$$

$$(1 \quad u_1 \quad \dots \quad u_{k-1}) (f_{k-w}(\mathbf{x}_i) + R);$$

where the last equality follows from $f_k(\mathbf{x}_i) + R = 0 \delta_k \in [K]$. Therefore,

$$\mathbb{A}(\mathbf{x}_i) = \sum_{k=1}^K (1 \quad u_1 \quad \dots \quad u_{k-1}) (f_k(\mathbf{x}_i) + R) = \sum_{k=1}^K (1 \quad u_1 \quad \dots \quad u_{k-1}) (f_{k-1}(\mathbf{x}_i) + R) + \mathbf{e}_{D-2}$$

$$= \mathbb{P}(\mathbf{x}_i) \mathbf{e}_{D-2};$$

where the last equality is due to (41). Thus letting the attention layer $\text{Attn}^{(3)} = f(\mathbf{V}_{k;w}; \mathbf{Q}_{k;w}; \mathbf{K}_{k;w})g_{(k;w) \in [K+1] \times \{0,1\}}$, we have

$$\begin{aligned} \mathbf{h}_i^{(3)} &= \text{Attn}(\mathbf{H}^{(2)})_i = \mathbf{h}_i + \frac{1}{N+1} \sum_{j=1}^{N+1} \sum_{k;w} \mathbf{Q}_{k;w} \mathbf{h}_j^{(2)}; \mathbf{K}_{k;w} \mathbf{h}_j^{(2)} \mathbf{V}_{k;w} \mathbf{h}_j^{(2)} \\ &= \mathbf{h}_i^{(2)} + \mathbb{P}(\mathbf{x}_i) \mathbf{e}_{D-2} \\ &= [\mathbb{P}(\mathbf{x}_i); f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); u_1; \dots; u_K; \mathbb{P}(\mathbf{x}_i); 1; t_i]; \end{aligned}$$

This is the desired result. \square

Now, we are ready to prove Theorem I.1.

Proof of Theorem I.1 As $\mathbb{P}(\cdot)$ is $(\epsilon=3; R; M; C)$ -approximable by sum of relus, we can invoke Proposition I.1 to show that there exists a single attention layer $\text{Attn}^{(1)}$ so that $\text{Attn}^{(1)}$ maps

$$\mathbf{h}_i \mapsto \mathbf{h}'_i = [\mathbf{x}_i; y_i; \mathbb{P}(\mathbf{x}_i); f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbb{E}_{\text{val}}(f_1); \dots; \mathbb{E}_{\text{val}}(f_K); 0; 1; t_i]; \quad i \geq [N+1];$$

for any input $\mathbf{H} = [\mathbf{h}_i]_i$ of the form described in Theorem I.1, and $\mathbb{E}_{\text{val}}(\cdot)$ is a functional such that $\max_k \mathbb{E}_{\text{val}}(f_k) - \mathbb{E}_{\text{val}}(f_k) = \epsilon$.

Next, by the proof of Proposition I.2, there exists $(\text{Attn}^{(1)}; \text{mlp}^{(2)}; \text{mlp}^{(3)})$ that maps

$$\mathbf{h}'_i \mapsto \mathbf{h}_i^{(3)} = [\mathbf{x}_i; y_i; \mathbb{P}(\mathbf{x}_i); f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \sum_{k=1}^K f_k(\mathbf{x}_i); 1; t_i]; \quad i \geq [N+1];$$

where $\mathbb{P}(\cdot) = (\mathbb{P}_1; \dots; \mathbb{P}_K) \in ([K])$ and $\mathbb{P}_k > 0$ only when $\mathbb{E}_{\text{val}}(f_k) = \min_{k' \neq k} \mathbb{E}_{\text{val}}(f_{k'}) + \epsilon$. Using the fact that $\max_k \mathbb{E}_{\text{val}}(f_k) - \mathbb{E}_{\text{val}}(f_k) = \epsilon$, we deduce that \mathbb{P} is supported on $f_k : \mathbb{E}_{\text{val}}(f_k) = \min_{k' \in [K]} \mathbb{E}_{\text{val}}(f_{k'}) + \epsilon$.

Therefore, $(\text{Attn}^{(1)}; \text{mlp}^{(1)}; \text{mlp}^{(2)}; \text{mlp}^{(3)})$ is the desired transformer, with

$$\max_{\epsilon \in [3]} M^{(\epsilon)} \leq (M+3)K; \quad \max_{\epsilon \in [3]} D^{(\epsilon)} \leq K^2 + K + 1;$$

and

$$\begin{aligned} \mathbb{J} \mathbb{J} \quad \max \quad & 3R + \frac{2NKC}{jD_{\text{val}}j} + 3K + 1; K + 3 + \epsilon^{-1}; KR + 2K + 2 \\ & 7KR + \frac{2NKC}{jD_{\text{val}}j} + \epsilon^{-1}; \end{aligned}$$

This completes the proof. \square

I.2 Proof of Theorem 11

We first restate Theorem 11 into the following version which provides additional size bounds for \mathcal{F} . For the simplicity of presentation, throughout this subsection and Appendix J, we denote $I_t = f_i : (\mathbf{x}_i; y_i) \in D_{\text{train}}g$, $I_v = f_i : (\mathbf{x}_i; y_i) \in D_{\text{val}}g$, $\mathbf{X}_{\text{train}} = [\mathbf{x}_i]_{i \in I_t}$ to be the input matrix corresponding to the training split only, and $N_{\text{train}} = jD_{\text{train}}j$, $N_{\text{val}} = jD_{\text{val}}j$.

Theorem I.2. For any sequence of regularizations $f_k g_{k \in [K]}$, 0 with $\epsilon := \max_k \frac{\epsilon + k}{k}$, $B_w > 0$, $\epsilon > 0$, and $\epsilon < B_w = 2$, suppose in input format (3) we have $D = (Kd)$. Then there exists an L -layer transformer TF with

$$L = d \log(B_w = (2^\epsilon))e + 4; \quad \max_{\epsilon \in [L]} M^{(\epsilon)} \leq 3K + 1; \quad \max_{\epsilon \in [L]} D^{(\epsilon)} \leq K^2 + K + 1;$$

$$\| \mathbf{w} \| \leq \frac{1}{\sqrt{K}} \left(\frac{1}{\sqrt{N_{\text{val}}}} + \frac{1}{\sqrt{N}} \right)^{-1} ; \quad R := \max \{ B_X B_W, B_Y \} ;$$

such that the following holds. On any input data $(D; \mathbf{x}_{N+1})$ such that the problem (ICRidge) is well-conditioned and has a bounded solution:

$$\min_{\mathbf{w}} (\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} = N_{\text{train}}) \quad \max_{\mathbf{w}} (\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} = N_{\text{train}}) \quad ; \quad \max_{k \in [K]} \mathbf{w}_{\text{ridge}}^k(D_{\text{train}}) \leq B_W = 2; \quad (42)$$

TF^0 approximately implements ridge selection: its prediction

$$\hat{y}_{N+1} = \text{read}_y(\text{TF}^0(\mathbf{H})) = \sum_{k=1}^K \mathbf{w}_k^\top \mathbf{x}_{N+1}; \quad \mathbf{w} = \sum_{k=1}^K \mathbf{w}_k$$

satisfies the following.

1. For each $k \in [K]$, $\mathbf{w}_k = \mathbf{w}_k(D_{\text{train}})$ approximates the ridge estimator $\mathbf{w}_{\text{ridge}}^k(D_{\text{train}})$, i.e. $\| \mathbf{w}_k - \mathbf{w}_{\text{ridge}}^k(D_{\text{train}}) \| \leq \epsilon$.
2. $\mathbf{w}_k \neq \mathbf{0}$ only if $\mathcal{L}_{\text{val}}(\mathbf{w}_k) \leq \min_{k' \in [K]} \mathcal{L}_{\text{val}}(\mathbf{w}_{k'}) + \epsilon$.

In particular, if we set $\epsilon = 2(B_X B_W + B_Y) B_X^{-1}$, then it holds that⁸

$$\text{dist}(\mathbf{w}; \text{conv} \{ \mathbf{w}_{\text{ridge}, \text{train}}^k \}) \leq \mathcal{L}_{\text{val}}(\mathbf{w}_{\text{ridge}, \text{train}}^k) \leq \min_{k' \in [K]} \mathcal{L}_{\text{val}}(\mathbf{w}_{\text{ridge}, \text{train}}^{k'}) + \epsilon;$$

where we denote $\mathbf{w}_{\text{ridge}, \text{train}}^k := \mathbf{w}_{\text{ridge}}^k(D_{\text{train}})$.

To prove Theorem I.2, we first show that, for the squared validation loss, there exists a 3-layer transformer that performs predictor selection based on the *exactly* evaluated $\mathcal{L}_{\text{val}}(f_k)$ for each $k \in [K]$. (Proof in Appendix I.2.1.)

Theorem I.3 (Square-loss version of Theorem I.1). *Consider the squared validation loss*

$$\mathcal{L}_{\text{val}}(f) := \frac{1}{2jD_{\text{val}}} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}_{\text{val}}} (f(\mathbf{x}_i) - y_i)^2;$$

Then there exists a 3-layer transformer TF with

$$\max_{\tau \in [3]} M^{(\tau)} \leq 2K + 2; \quad \max_{\tau \in [3]} D^{(\tau)} \leq K^2 + K + 1; \quad \| \mathbf{w} \| \leq \frac{1}{\sqrt{K}} \left(\frac{1}{\sqrt{jD_{\text{val}}}} + \frac{1}{\sqrt{N}} \right)^{-1};$$

such that for any input \mathbf{H} that takes form

$$\mathbf{h}_i = [\mathbf{x}_i; y_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbf{0}_K; 1; t_i];$$

where TF outputs $\mathbf{h}_{N+1} = [\mathbf{x}_{N+1}; \hat{y}(\mathbf{x}_{N+1}); 1; 0]$, where the predictor $\hat{y}: \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex combination of $f_k: \mathbb{R}^d \rightarrow \mathbb{R}$ with $\mathcal{L}_{\text{val}}(f_k) \leq \min_{k' \in [K]} \mathcal{L}_{\text{val}}(f_{k'}) + \epsilon$. As a corollary, for any convex risk $L: (\mathbb{R}^d \rightarrow \mathbb{R}) \rightarrow \mathbb{R}$, \hat{y} satisfies

$$L(\hat{y}) \leq \min_{k' \in [K]} L(f_{k'}) + \max_{k \in [K]} \mathcal{L}_{\text{val}}(f_k) \leq L(f_k) + \epsilon;$$

Proof of Theorem I.2 First, by the proof⁹ of Theorem 4 and Proposition B.6, for each $k \in [K]$, there exists a $T = L + 3$ layer transformer $\text{TF}^{(1:T)}$ such that $\text{TF}^{(1:T)}$ maps

$$\mathbf{h}_i \mapsto \mathbf{h}_i^{(T)} = [\mathbf{x}_i; y_i; \mathbf{w}_1^\top \mathbf{x}_i; \dots; \mathbf{w}_K^\top \mathbf{x}_i; \mathbf{0}_K; 1; t_i];$$

⁸This is because $\mathcal{L}_{\text{val}}(\mathbf{w})$ is $(B_X B_W + B_Y) B_X^{-1}$ -Lipschitz w.r.t. $\mathbf{w} \in \mathcal{B}_2(B_W)$.

⁹Technically, an adapted version where the underlying ICGD mechanism operates on the training split (with $t_i = 1$) with size N_{train} instead of on all N training examples, which only changes $\| \mathbf{w} \|$ by at most a constant factor, and does not change the number of layers and heads.

so that if (42) holds, we have $\|\mathbf{w}_k - \mathbf{w}_{\text{ridge}}^k\|_2 \leq \epsilon$ and $\mathbf{w}_k \in B_2(B_w)$.

Next, by Theorem I.3, there exists a 3-layer transformer $\text{TF}^{(T+1; T+3)}$ that outputs

$$\mathbf{h}_{N+1}^{(T+3)} = [\mathbf{x}_{N+1}; \mathbf{h}; \mathbf{x}_{N+1}; \mathbf{i}; \mathbf{1}; t_i];$$

where $\mathbf{w} = \prod_{k=1}^K \mathbf{w}_k$, $\mathbf{w} = (\mathbf{w}_1; \dots; \mathbf{w}_K) \in ([K])$ so that

$$k > 0 \text{ only if } \mathbb{E}_{\text{val}}(\mathbf{w}_k) = \min_{k' \in [K]} \mathbb{E}_{\text{val}}(\mathbf{w}_{k'}) + \epsilon.$$

This is the desired result. \square

I.2.1 Proof of Theorem I.3

Similar to the proof of Proposition 10, Theorem I.3 is a direct corollary by combining Proposition I.3 with Proposition I.2.

Proposition I.3 (Evaluation layer for the squared loss). *There exists an attention layer TF with $2K$ heads and $\mathbb{J} \mathbb{J} = 3R + 2NK = jD_{\text{val}}$ such that TF maps*

$$\begin{aligned} \mathbf{h}_i &= [\mathbf{x}_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbf{0}_K; \mathbf{1}; t_i] \\ \mathbf{h}'_i &= [\mathbf{x}_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbb{E}_{\text{val}}(f_1); \dots; \mathbb{E}_{\text{val}}(f_K); \mathbf{1}; t_i]; \quad i \in [N+1]; \end{aligned}$$

Proof of Proposition I.3. For every $k \in [K]$, we define matrices $\mathbf{Q}_{m;k}, \mathbf{K}_{m;k}, \mathbf{V}_{m;k} \in \mathbb{R}^{D \times D}$ such that for all $i, j \in [N+1]$,

$$\begin{aligned} \mathbf{Q}_{k,0} \mathbf{h}_i &= \begin{bmatrix} 2 & 1 & 3 \\ 6 & 17 & 5 \\ 0 & 0 & 0 \end{bmatrix}; \quad \mathbf{Q}_{k,1} \mathbf{h}_i = \begin{bmatrix} 2 & 1 & 3 \\ 6 & 17 & 5 \\ 0 & 0 & 0 \end{bmatrix}; \quad \mathbf{K}_{k,0} \mathbf{h}_j = \mathbf{K}_{k,1} \mathbf{h}_j = \begin{bmatrix} 2 & f_k(\mathbf{x}_j) & 3 \\ 6 & y_j & 5 \\ 0 & 0 & 0 \end{bmatrix}; \\ \mathbf{V}_{k,0} \mathbf{h}_j &= \mathbf{V}_{k,1} \mathbf{h}_j = \frac{(N+1)}{2jD_{\text{val}}} (f_k(\mathbf{x}_j) \quad y_j) \mathbf{e}_{D-(K-k)-3}. \end{aligned}$$

As the input has structure $\mathbf{h}_i = [\mathbf{x}_i; y_i; \dots; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbf{0}_{K+1}; \mathbf{1}; t_i]$, these matrices indeed exist, and further it is straightforward to check that they have norm bounds

$$\max_{k \in [K]; w \in \{0,1\}} \|\mathbf{Q}_{k,w}\|_{\text{op}} \leq 3; \quad \max_{k \in [K]; w \in \{0,1\}} \|\mathbf{K}_{k,w}\|_{\text{op}} \leq 1 + R; \quad \prod_{k \in [K]; w \in \{0,1\}} \|\mathbf{V}_{k,w}\|_{\text{op}} \leq \frac{K(N+1)}{jD_{\text{val}}}.$$

Now, for every $i, j \in [N+1]$, we have

$$\begin{aligned} & \prod_{w \in \{0,1\}} (\mathbf{h} \mathbf{Q}_{k,w} \mathbf{h}_i; \mathbf{K}_{k,w} \mathbf{h}_j) \mathbf{V}_{k,w} \mathbf{h}_j \\ &= [(f_k(\mathbf{x}_j) \quad y_j \quad 2R(1+t_j)) \quad (y_j \quad f_k(\mathbf{x}_j) \quad 2R(1+t_j))] \frac{(N+1)}{2jD_{\text{val}}} (f_k(\mathbf{x}_j) \quad y_j) \mathbf{e}_{D-(K-k)-3} \\ &= 1ft_j = 1g [(f_k(\mathbf{x}_j) \quad y_j) \quad (y_j \quad f_k(\mathbf{x}_j))] \frac{(N+1)}{2jD_{\text{val}}} (f_k(\mathbf{x}_j) \quad y_j) \mathbf{e}_{D-(K-k)-3} \\ &= 1ft_j = 1g \frac{(N+1)}{2jD_{\text{val}}} (f_k(\mathbf{x}_j) \quad y_j)^2 \mathbf{e}_{D-(K-k)-3}; \end{aligned}$$

where the second equality follows from the bound $|f_k(\mathbf{x}_j) - y_j| \leq 2R$, so that the relus equals 0 if $t_j = 0$. Thus letting the attention layer $\text{Attn} = f(\mathbf{V}_{k,w}; \mathbf{Q}_{k,w}; \mathbf{K}_{k,w})_{(k,w) \in [K] \times \{0,1\}}$, we have

$$\hat{\mathbf{h}}_i = [\text{Attn}(\mathbf{H})]_i = \mathbf{h}_i + \frac{1}{N+1} \prod_{j=1}^{N+1} \prod_{k,w} (\mathbf{h} \mathbf{Q}_{k,w} \mathbf{h}_i; \mathbf{K}_{k,w} \mathbf{h}_j) \mathbf{V}_{k,w} \mathbf{h}_j$$

$$\begin{aligned}
&= \mathbf{h}_i + \frac{1}{2^j D_{\text{val}}^j} \sum_{j=1}^{K+1} \sum_{k=1}^K (f_k(\mathbf{x}_j) - y_j)^2 \mathbf{1}_{f_j} = \mathbf{1}_{g_{D-(K-k)-3}} \\
&= \mathbf{h}_i + \frac{1}{2^j D_{\text{val}}^j} \sum_{(\mathbf{x}_j, y_j) \in \mathcal{D}_{\text{val}}} (f_k(\mathbf{x}_j) - y_j)^2 \mathbf{A} \mathbf{e}_{D-(K-k)-3} \\
&= \mathbf{h}_i + \sum_{k=1}^K \mathbf{e}_{\text{val}}(f_k) \mathbf{e}_{D-(K-k)-3} \\
&= [\mathbf{x}_i; y_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbf{0}_{K+1}; 1; t_i] + [\mathbf{0}_{D-K-3}; \mathbf{e}_{\text{val}}(f_1); \dots; \mathbf{e}_{\text{val}}(f_K); 0; 0; 0] \\
&= [\mathbf{x}_i; y_i; f_1(\mathbf{x}_i); \dots; f_K(\mathbf{x}_i); \mathbf{e}_{\text{val}}(f_1); \dots; \mathbf{e}_{\text{val}}(f_K); 0; 1; t_i]; \quad i \geq [N+1];
\end{aligned}$$

This is the desired result. \square

I.3 Proofs for Section 4.2

I.3.1 Proof of Lemma 13

It is straightforward to check that the binary type check $\text{binary} : \mathbb{R} \rightarrow \mathbb{R}$ can be expressed as a linear combination of 6 relu's (recalling $\text{relu}(x) = \text{ReLU}(x)$):

$$\begin{aligned}
\text{binary}(y) &= \frac{y+1}{2} - \frac{y}{2} + \frac{y-1}{2} + \frac{y}{2} \left(\frac{1-y}{2} \right) - \frac{y-1}{2} + \frac{y}{2} \left(\frac{1+y}{2} \right) \\
&= \sum_{m=1}^6 a_m \text{relu}(b_m y + c_m);
\end{aligned}$$

with $\sum_{m=1}^6 |a_m| = 8$, $\max_m \max_j |b_m|, |c_m| \leq 2$. We can thus construct an attention layer $f : (\mathbf{Q}_m; \mathbf{K}_m; \mathbf{V}_m) \rightarrow \mathbb{R}^6$ with 6 heads such that

$$\mathbf{Q}_m \mathbf{h}_i = [b_m; c_m; \mathbf{0}_{D-2}]; \quad \mathbf{K}_m \mathbf{h}_j = [y_j; 1; \mathbf{0}_{D-2}]; \quad \mathbf{V}_m \mathbf{h}_j = \frac{N+1}{N} a_m t_j; \mathbf{0}_{D-1};$$

which gives that for every $i \geq [N+1]$,

$$\begin{aligned}
&\sum_{m=1}^6 \frac{1}{N+1} \sum_{j \in [N+1]} (t \mathbf{Q}_m \mathbf{h}_i; \mathbf{K}_m \mathbf{h}_j) [\mathbf{V}_m \mathbf{h}_j]_1 \\
&= \sum_{m=1}^6 \frac{1}{N} \sum_{j=1}^N (b_m y_j + c_m) a_m = \frac{1}{N} \sum_{j=1}^N \text{binary}(D);
\end{aligned}$$

Further, we have $\sum_{m=1}^6 |a_m| = 8 = O(1)$. This is the desired result. \square

By composing the above attention layer with one additional layer (with 2 heads) that implement the following function

$$(2(t-1/2)) \wedge (2(t-1));$$

on the output $\text{binary}(D)$, we directly obtain the following corollary.

Corollary I.1 (Thresholded binary test). *There exists a two-layer attention-only transformer with $\max_{i \in [2]} M^{(i)} \leq 6$ and $\sum_{i=1}^2 O(1) = O(1)$ that exactly implements the thresholded binary test*

$$\text{binary}_{\text{thres}}(D) := \begin{cases} 1; & \text{if } \text{binary}(D) \geq 1; \\ 0; & \text{if } \text{binary}(D) \leq \frac{1}{2}; \\ \text{linear interpolation}; & \text{o.w.} \end{cases} \quad (43)$$

at every token $i \geq [N+1]$, where we recall the definition of binary in Lemma 13.

I.3.2 Formal statement and proof of Proposition 14

We say a distribution P_y on \mathbb{R} is $(C; \epsilon_0)$ -not-concentrated around $f_0; 1g$ if

$$P_y([-\epsilon_0; \epsilon_0] \cap [1 - \epsilon_0; 1 + \epsilon_0]) \geq C$$

for all $\epsilon \geq (0; \epsilon_0]$. A sufficient condition is that the density p_y is upper bounded by C within $[-\epsilon_0; \epsilon_0] \cap [1 - \epsilon_0; 1 + \epsilon_0]$.

Throughout this section, let $\sigma(t) := (1 + e^{-t})^{-1}$ denote the sigmoid activation, and let \mathbf{w}_{\log} denote the solution to the in-context logistic regression problem, i.e. (ICGLM) with $g(\cdot) = \sigma(\cdot)$.

Proposition I.4 (Adaptive regression or classification; Formal version of Proposition 14). *For any $B_w > 0$, $\epsilon \in (0; 1]$, $B_x B_w = 10$, $0 < \epsilon_0 < 1$ with $\epsilon := \epsilon_0$, and any $(C; \epsilon_0)$, there exists a L -layer attention-only transformer with*

$$L \leq O\left(\log \frac{B_x B_w}{\epsilon}\right); \max_{i \in [L]} M^{(i)} \leq O\left(1 + \frac{B_x^4}{2} \epsilon^{-2}\right); \|\mathbf{w}_{\log}\| \leq O\left(R + \frac{1}{\epsilon} + \frac{1}{\epsilon}\right)$$

(with $R := \max\{B_x B_w; B_y; 1\}$, and ϵ depending only on $(C; \epsilon_0)$) such that the following holds. Suppose the input format is (3) with dimension $D \geq 3d + 4$.

On any classification instance $(D; \mathbf{x}_{N+1})$ (such that $f_{y_i} g_{i \in [N]}(f_0; 1g)$ that is well-conditioned for logistic regression in the sense of (35), it outputs \hat{y}_{N+1} that ϵ -approximates the prediction of in-context logistic regression:

$$\hat{y}_{N+1} = \sigma(\langle \mathbf{h}_{\mathbf{x}_{N+1}}; \mathbf{w}_{\log} \rangle)$$

On the contrary, for regression problems, i.e. any in-context distribution P whose marginal P_y is $(C; \epsilon_0)$ -not-concentrated around $f_0; 1g$, with probability at least $1 - \exp(-cN)$ over D (where $c > 0$ depends only on $(C; \epsilon_0)$), \hat{y}_{N+1} ϵ -approximates the prediction of in-context least squares if the data is well-conditioned:

$$\hat{y}_{N+1} = \langle \mathbf{h}_{\mathbf{x}_{N+1}}; \mathbf{w}_{LS} \rangle \text{ whenever } D \text{ satisfies (5) with } \epsilon = 0;$$

where \mathbf{w}_{LS} denotes the in-context least squares estimator, i.e. (ICRidge) with $\epsilon = 0$.

Proof. The result follows by combining the binary test in Corollary I.1 with Theorem 4 and Theorem G.1. By those results, there exists three attention-only transformers $\mathcal{L}_{LS}; \mathcal{L}_{\log}; \mathcal{L}_{bin}$, with (below $L_g; C_g = (1)$ for $g = \sigma(\cdot)$)

$$\begin{aligned} L_{LS} &\leq O\left(\log \frac{B_x B_w}{\epsilon}\right); \max_{i \in [L_{LS}]} M_{LS}^{(i)} \leq 3; \|\mathbf{w}_{LS}\| \leq O\left(R + \frac{1}{\epsilon}\right); \\ L_{\log} &\leq O\left(\log \frac{L_g B_x B_w}{\epsilon}\right); \max_{i \in [L_{\log}]} M_{\log}^{(i)} \leq C_g^2 \left(1 + \frac{L_g^2 B_x^4}{2} \epsilon^{-2}\right); \|\mathbf{w}_{\log}\| \leq O\left(R + \frac{C_g}{\epsilon}\right); \\ L_{bin} &= 2; \max_{i \in [2]} M_{bin}^{(i)} \leq 6; \|\mathbf{w}_{bin}\| \leq O(1/\epsilon); \end{aligned}$$

that outputs prediction \hat{y}_{N+1}^{LS} , \hat{y}_{N+1}^{\log} (at the $(N + 1)$ -th token) and $\text{binary}_{thres}(D)$ (at every token) respectively, which satisfy

$$\begin{aligned} \hat{y}_{N+1}^{\log} &= \sigma(\langle \mathbf{h}_{\mathbf{x}_{N+1}}; \mathbf{w}_{\log} \rangle) \\ \hat{y}_{N+1}^{LS} &= \langle \mathbf{h}_{\mathbf{x}_{N+1}}; \mathbf{w}_{LS} \rangle \end{aligned}$$

when the corresponding well-conditionednesses are satisfied. In particular, we can make \mathbf{w}_{\log} well-defined on non-binary data, by multiplying $\text{binary}_{thres}(D)$ onto the \mathbf{x}_i 's (which can be implemented by slightly modifying \mathcal{L}_{\log} without changing the order of the number of layers, heads, and norms) so that $\mathbf{w}_{\log} = \mathbf{0}$ on any data where $\text{binary}_{thres}(D) = 0$.

By joining \mathcal{L}_{LS} and \mathcal{L}_{\log} using Proposition B.6, concatenating with \mathcal{L}_{bin} before, and concatenating with one additional attention layer with 2 heads after to implement

$$\text{binary}_{thres}(D) \hat{y}_{N+1}^{\log} + 1 - \text{binary}_{thres}(D) \hat{y}_{N+1}^{LS} \quad (44)$$

we obtain a single transformer with

$$L = O\left(\log \frac{B_x B_w}{\epsilon}\right); \max_{\epsilon \in [L]} M^{(\epsilon)} = O\left(1 + \frac{B_x^4}{2} \epsilon^{-2}\right); \|\mathbb{J}_{LS}\| = O\left(R + \frac{1}{\epsilon} + \frac{1}{\epsilon}\right);$$

which outputs (44) as its prediction (at the location for y_{N+1}).

It remains to show that (44) reduces to either one of y_{N+1}^{\log} or y_{N+1}^{LS} . When the data are binary ($y_i \in \{0, 1\}$), we have $\text{binary}(D) = 1$ and $\text{binary}_{\text{thres}}(D) = 1$, in which case (44) becomes exactly y_{N+1}^{\log} . By contrast, when data is sampled from a distribution that is $(C; \epsilon)$ -not-concentrated around $f_0; 1g$, we have for any fixed $\epsilon \in [L]$ that, letting $B_\epsilon := [\epsilon; 1] [1 - \epsilon; 1 + \epsilon]$ and $p_\epsilon := P_y(B_\epsilon) = C \epsilon^{-\frac{1}{4}}$, by Hoeffding's inequality,

$$\begin{aligned} P\left(\text{binary}_{\text{thres}}(D) \neq 0\right) &= P\left(\text{binary}_{\text{thres}}(D) \leq \frac{1}{2}\right) = P\left(\frac{1}{N} \sum_{i=1}^N 1fy_i \geq B_\epsilon g \leq \frac{1}{2}\right) \\ &\leq \exp(-c(1-2p_\epsilon)^2 N) = \exp(-c'N); \end{aligned}$$

where $c' > 0$ is an absolute constant. On the event $\text{binary}_{\text{thres}}(D) = 0$ (which happens with probability at least $1 - \exp(-c'N)$), (44) becomes exactly y_{N+1}^{LS} . This finishes the proof. \square

I.4 Linear correlation test and application

In this section, we give another instantiation of the pre-ICL testing mechanism by showing that the transformer can implement a *linear correlation test* that tests whether the correlation vector $\mathbb{E}[xy]$ has a large norm. We then use this test to construct a transformer to perform ‘‘confident linear regression’’, i.e. output a prediction from linear regression only when the signal-to-noise ratio is high.

For any fixed parameters $\min(B_w^2) > 0$, consider the linear correlation test over data D defined as

$$\begin{aligned} \text{lin}(D) &:= \frac{1}{\sum_{\min(B_w^2)=2}^h} \begin{cases} k_2^2 & (\min(B_w^2)=4)^2 \\ k_2^2 & (3 \min(B_w^2)=4)^2 \\ & \vdots \end{cases} \\ &\leq 0; \\ &= 1; \\ &\text{linear interpolation}; \text{ o.w.}; \end{aligned} \tag{45}$$

$$\text{where } \mathbb{K} = \mathbf{T}(D) := \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i y_i;$$

Recall that $\text{ReLU}(\cdot) = \max(\cdot, 0)$ above denotes the relu activation.

We show that lin can be exactly implemented by a 3-layer transformer.

Lemma I.1 (Expressing lin by transformer). *There exists a 3-layer attention-only transformer TF with at most 2 heads per layer and $\|\mathbb{J}\| = O\left(1 + \frac{2}{\min(B_w^2)}\right)$ such that on input sequence \mathbf{H} of the form (3) with $D = 2d + 4$, the transformer exactly implements lin : it outputs \mathbf{H} such that $\mathbf{H}_i = [\mathbf{x}_i; y_i t_i; \dots; \text{lin}(D); 1]$ for all $i \in [N + 1]$.*

Proof. We begin by noting the following basic facts:

- Identity function can be implemented exactly by two ReLUs: $t = \text{ReLU}(t) + \text{ReLU}(-t)$.
- Squared ℓ_2 norm can be implemented exactly by a single attention head (assuming every input \mathbf{h}_i contains the same vector \mathbf{g}): $\|\mathbf{g}\|_2^2 = \text{ReLU}(\mathbf{g}; \mathbf{g})$.

We construct the transformer as follows.

Layer 1: Use 2 heads to implement $\mathbb{K} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i y_i$, where $\mathbf{V}_{\{1,2\}}^{(1)} \mathbf{h}_j = [\mathbf{x}_j; \mathbf{0}_{D-d}]$, $\mathbf{Q}_{\{1,2\}}^{(1)} \mathbf{h}_i = [\frac{N+1}{N}; \mathbf{0}_{D-1}]$, and $\mathbf{K}_{\{1,2\}}^{(1)} \mathbf{h}_j = [y_j t_j; \mathbf{0}_{D-1}] = [y_j 1fj < N + 1g; \mathbf{0}_{D-1}]$ (where we recall $t_j = 1fj < N + 1g$ and note that $y_j t_j$ corresponds exactly to the location for y_j in \mathbf{H} , cf. (3)).

On the above event, we have

$$\mathbb{E}_2 = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i (\mathbf{h} \mathbf{x}_i; \mathbf{w}_P^? i + \mathbf{z}_i) = \mathbf{b} \mathbf{w}_P^? + \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_i :$$

Therefore, in case 1, we have

$$\mathbb{E}_2 \mathbf{b} \mathbf{w}_P^? \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_i \geq 0.9 \min \{k \mathbf{w}_P^? k_2, \frac{\min B_W^?}{8}\} \frac{3 \min B_W^?}{4}.$$

In case 2, we have

$$\mathbb{E}_2 \mathbf{b} \mathbf{w}_P^? + \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{z}_i \leq \max \left\{ \frac{\min B_W^?}{10 \max}, \frac{\min B_W^?}{8}, \frac{\min B_W^?}{4} \right\}.$$

The proof is finished by recalling the definition of lin in (45), so that $\text{lin}(D) = 1$ if $k \mathbf{w}_P^? k_2 \geq \frac{3 \min B_W^?}{4}$, and $\text{lin}(D) = 0$ if $k \mathbf{w}_P^? k_2 < \frac{\min B_W^?}{4}$. \square

Application: Confident linear regression By directly composing the linear correlation test in Lemma I.1 with the transformer construction in Corollary 5 (using an argument similar as the proof of Proposition I.4), and using the power of the linear correlation test Proposition I.5, we immediately obtain the following result, which outputs a prediction from (approximately) least squares if $\text{lin}(D) = 1$, and abstains from predicting if $\text{lin}(D) = 0$. This can be viewed as a form of “confident linear regression”, where the model predicts only if it thinks the linear signal is strong enough.

Proposition I.6 (Confident linear regression). *For any $B_W > 0, 0 < B_W^? \leq \frac{B_W}{10}$, $0 < \frac{B_W^?}{B_W} \leq 1$, with $\frac{B_W^?}{B_W} := \frac{\min B_W^?}{\max B_W^?}$, there exists a L -layer attention-only transformer with*

$$L = O \left(\log \frac{B_X B_W}{B_W^?} \right); \max_{i \in [L]} M^{(i)} = O(1); \mathbb{E} \left[\sum_{i=1}^L R + \frac{1}{\min(B_W^?)^2} \right]$$

(with $R := \max \{f(B_X B_W; B_Y; 1)g\}$) such that the following holds. Let $N \geq \max \{K^4, \frac{\max K^2}{(B_W^?)^2} \frac{2}{\min} g\} d$. Suppose the input format is (3) with dimension $D = 2d + 4$. Let ICL instance $(D; \mathbf{x}_{N+1})$ be drawn from any distribution \mathbb{P} satisfying Assumption D. Then the transformer outputs a 2-dimensional prediction (within the test token \mathbf{h}_{N+1})

$$(\mathbf{y}_{N+1}; \mathbf{b}) \in \mathbb{R}^2 \times \{0, 1\}^g$$

such that the following holds:

1. If $k \mathbf{w}_P^? k_2 \geq \frac{3 \min B_W^?}{4}$, then with probability at least $1 - \frac{1}{N}$ over D , we have $\mathbf{y}_{N+1} = \mathbf{h}_{N+1} \mathbf{w}_{LS}; \mathbf{x}_{N+1} \mathbf{j}$; and $\mathbf{b} = 1$ if D is in addition well-conditioned for least squares (in the sense of (5) with $\gamma = 0$).
2. If $k \mathbf{w}_P^? k_2 < \frac{3 \min B_W^?}{4}$, then with probability at least $1 - \frac{1}{N}$ over D , we have $\mathbf{y}_{N+1} = 0$ and $\mathbf{b} = 0$.

J Proof of Theorem 12: Noisy linear model with mixed noise levels

For each fixed $k \geq [K]$, we consider the following data generating model \mathbb{P}_k , where we first sample $\mathbf{P} = \mathbb{P}_{\mathbf{w}_?; k}$ from $\mathbf{w}_? \sim \mathcal{N}(\mathbf{0}; \mathbf{I}_d)$, and then sample data $f(\mathbf{x}_i; y_i)_{i \in [N+1]} \stackrel{\text{iid}}{\sim} \mathbb{P}_{\mathbf{w}_?; k}$ as

$$\mathbb{P}_{\mathbf{w}_?; k} : \mathbf{x}_i \sim \mathcal{N}(\mathbf{0}; \mathbf{I}_d); y_i = \mathbf{h} \mathbf{x}_i; \mathbf{w}_? i + \epsilon_i; \epsilon_i \sim \mathcal{N}(0; \frac{2}{k});$$

Also, recall that the Bayes optimal estimator on \mathbb{P}_k is given by $\mathbf{y}_{N+1}^{\text{Bayes}} = \mathbf{w}_{\text{ridge}}^k(D); \mathbf{x}_{N+1}$ with ridge $\lambda = \frac{2}{k} d = N$, and the Bayes risk on \mathbb{P}_k is given by

$$\text{BayesRisk}_k := \inf_{\mathcal{A}} \mathbb{E}_k \frac{1}{2} (A(D)(\mathbf{x}_{N+1}) - y_{N+1})^2 = \mathbb{E}_k \frac{1}{2} \mathbf{y}_{N+1}^{\text{Bayes}} - y_{N+1}^2 :$$

Recall that in Section 4.1.1, we consider a mixture law \mathbb{P} that generates data from \mathbb{P}_k with $k \in [K]$. It is clear that we have (pushing $\inf_{\mathcal{A}}$ into $E_{k \sim \mathbb{P}}$ does not increase the value) we have

$$\text{BayesRisk} \leq E_{k \sim \mathbb{P}} [\text{BayesRisk}_k];$$

i.e., the Bayes risk can only be greater if we consider a mixture of models. In other words, if a transformer can achieve near-Bayes ICL on each meta-task \mathbb{P}_k , then it can perform near-Bayes ICL on any meta-task which is a mixture of \mathbb{P}_k with $k \in [K]$. Therefore, to prove Theorem 12, it suffices to show the following (strengthened) result.

Theorem J.1 (Formal version of Theorem 12). *Suppose that $N \geq 0.1d$ and we write $\gamma_{\max} = \max_k f_{k;1}g$; $\gamma_{\min} = \min_k f_{k;1}g$. Suppose in input format (3) we have $D \geq (Kd)$. Then there exists a transformer with*

$$L \leq O\left(\frac{1}{\gamma_{\min}} \log(N - \gamma_{\min})\right); \quad \max_{\ell \in [L]} M^{(\ell)} \leq O(K); \quad \max_{\ell \in [L]} D^{(\ell)} \leq O(K^2);$$

$$\| \mathbb{Y}_{N+1} - \hat{\mathbb{Y}}_{N+1} \| \leq O\left(\frac{1}{\gamma_{\max}} K d \log(N)\right);$$

such that for any $k \in [K]$, it holds that

$$E_k \frac{1}{2} (\mathbb{Y}_{N+1} - \hat{\mathbb{Y}}_{N+1})^2 \leq \text{BayesRisk}_k + \Theta\left(\frac{2}{\gamma_{\min}} \frac{\log K}{N}\right)^{1=3}$$

if we choose $N_{\text{val}} := jD_{\text{val}}j \geq N^{2=3}[\log K]^{1=3}$.

The core of the proof of Theorem J.1 is to show that any estimator $\hat{\mathbb{Y}}$ that achieves small validation loss \mathbb{E}_{val} must achieve small population loss.

Throughout the rest of this section, recall that we define $N_{\text{train}} = jD_{\text{train}}j$; $N_{\text{val}} = jD_{\text{val}}j$, $l_{\text{train}} = fi : (\mathbf{x}_i; y_i) \in D_{\text{train}}g$, $l_{\text{val}} = fi : (\mathbf{x}_i; y_i) \in D_{\text{val}}g$, and $\mathbf{X}_{\text{train}} = [\mathbf{x}_i]_{i \in \mathcal{I}_{\text{train}}}$.

J.1 Proof of Theorem J.1

Fix parameters $\epsilon, \delta, \eta > 0$ and a large universal constant C_0 . Let us set

$$r = \max\left\{\frac{1}{\epsilon}, \frac{1}{\delta}\right\}; \quad \rho = \frac{1}{dN_{\text{train}}}; \quad \kappa = 25;$$

$$B_w^? = 1 + C_0 \frac{r \log(N)}{d}; \quad B_w = C_0(B_w^? + \gamma_{\max});$$

$$B_x = C_0 \frac{\rho}{d \log(N)}; \quad B_y = C_0(B_x^? + \gamma_{\max}) \frac{\rho}{\log(N)};$$

Then, we define good events similarly to the proof of Corollary 6 (Appendix F.4):

$$E = \{k \leq \kappa, B_w^? \leq \kappa, \rho \leq \frac{1}{N}g\};$$

$$E_w = \{f \leq \min(\mathbf{X}_{\text{train}}^T \mathbf{X}_{\text{train}} = N_{\text{train}}) \max(\mathbf{X}_{\text{train}}^T \mathbf{X}_{\text{train}} = N_{\text{train}}) \leq g\};$$

$$E_{b;\text{train}} = \{f \leq \delta(\mathbf{x}_i; y_i) \leq D_{\text{train}}; k \leq \kappa, B_x; j y_i j \leq B_y g\};$$

$$E_{b;\text{val}} = \{f \leq \delta(\mathbf{x}_i; y_i) \leq D_{\text{val}}; k \leq \kappa, B_x; j y_i j \leq B_y g\};$$

$$E_{b;N+1} = \{f \leq \delta(\mathbf{x}_{N+1}; y_{N+1}) \leq D_{\text{val}}; k \leq \kappa, B_x; j y_{N+1} j \leq B_y g\};$$

For the good event $E := E \cap E_w \cap E_{b;\text{train}} \cap E_{b;\text{test}} \cap E_{b;N+1}$, we can show that $P(E^c) \leq O(N^{-10})$. Further, by the proof of Lemma F.1 (see e.g. (34)), we know that $\max_{k \in [K]} \mathbf{w}_{\text{ridge}}^k(D_{\text{train}}) \leq B_w = 2$ holds under the good event E .

For the ridge $\mathbf{w}_k = \frac{d}{N_{\text{train}}}$ and parameters $(\epsilon; \delta; \eta)$, we consider the transformer constructed in Theorem I.2, with a clipped prediction $\hat{\mathbb{Y}}_{N+1} = \text{clip}_{\gamma}(\text{TF}(\mathbf{H}))$.

In the following, we upper bound the quantity $E_k (\mathbb{Y}_{N+1} - \hat{\mathbb{Y}}_{N+1})^2$ for any fixed k . Similar to the proof of Corollary 6 (Appendix F.4), we decompose

$$E_k (\mathbb{Y}_{N+1} - \hat{\mathbb{Y}}_{N+1})^2 = E_k \mathbb{1}_{E^c} (\mathbb{Y}_{N+1} - \hat{\mathbb{Y}}_{N+1})^2 + E_k \mathbb{1}_E (\mathbb{Y}_{N+1} - \hat{\mathbb{Y}}_{N+1})^2;$$

and we analyze these two parts separately.

Part I. Recall that by our construction, when E holds, we have $\mathbf{y}_{N+1} = \text{clip}_{B_y}(h(\mathbf{w}; \mathbf{x}_{N+1}))$ and the statements of Theorem I.2 hold for \mathbf{w} . Thus, we have

$$\begin{aligned} E_k [1fEg(\mathbf{y}_{N+1} - y_{N+1})^2] &= E_k [1fEg(\text{clip}_{B_y}(h(\mathbf{x}_{N+1}; \mathbf{w})) - y_{N+1})^2] \\ &= E_k [1fEg(h(\mathbf{x}_{N+1}; \mathbf{w}) - y_{N+1})^2] : \end{aligned}$$

Let us consider the following risk functional

$$L_{\text{val}, \mathbf{w}_?}(\mathbf{w}) = E_{(\mathbf{x}; y) \sim P_{\mathbf{w}_?}} \frac{1}{2} (h(\mathbf{w}; \mathbf{x}) - y)^2 = \frac{1}{2} k \mathbf{w} - \mathbf{w}_? k_2^2 + \frac{2}{k} :$$

Then, under the good event $E_0 := E \setminus E_W \setminus E_{b, \text{train}} \setminus E_{b, \text{test}}$ of $(\mathbf{w}_?; D)$,

$$\begin{aligned} E_k [1fEg(h(\mathbf{x}_{N+1}; \mathbf{w}) - y_{N+1})^2 \mid \mathbf{w}_?; D] &= E_k [1fEg(h(\mathbf{x}_{N+1}; \mathbf{w}(D)) - y_{N+1})^2 \mid \mathbf{w}_?; D] \\ &= E_k [(h(\mathbf{x}_{N+1}; \mathbf{w}(D)) - y_{N+1})^2 \mid \mathbf{w}_?; D] \\ &= E_{(\mathbf{x}; y) \sim P_{\mathbf{w}_?}} (h(\mathbf{x}_{N+1}; \mathbf{w}(D)) - y_{N+1})^2 \\ &= L_{\text{val}, \mathbf{w}_?}(\mathbf{w}(D)) : \end{aligned}$$

By our construction, under the good event E_0 , we have

$$L_{\text{val}, \mathbf{w}_?}(\mathbf{w}(D)) = L_{\text{val}, \mathbf{w}_?}(\mathbf{w}_k(D_{\text{train}})) + \max_{l \in [K]} \mathbb{E}_{\text{val}}(\mathbf{w}_l(D_{\text{train}})) - L_{\text{val}, \mathbf{w}_?}(\mathbf{w}_l(D_{\text{train}})) + ;$$

where $\mathbb{E}_{\text{val}}(\mathbf{w}_l(D_{\text{train}})) = \mathbf{w}_{\text{ridge}}^l(D_{\text{train}})$ " for each $l \in [K]$. Clearly,

$$\begin{aligned} 2E_k [1fE_0g L_{\text{val}, \mathbf{w}_?}(\mathbf{w}_k(D_{\text{train}}))] &= E_k [1fE_0g k \mathbf{w}_k(D_{\text{train}}) - \mathbf{w}_? k_2^2 + \frac{2}{k}] \\ E_k [1fE_0g \mathbf{w}_{\text{ridge}}^k(D_{\text{train}}) - \mathbf{w}_? \frac{2}{2} + 2] &= \mathbf{w}_{\text{ridge}}^k(D_{\text{train}}) - \mathbf{w}_? \frac{2}{2} + \frac{2}{k} \\ E_k [\mathbf{w}_{\text{ridge}}^k(D_{\text{train}}) - \mathbf{w}_? \frac{2}{2} + 2] &= \mathbf{w}_{\text{ridge}}^k(D_{\text{train}}) - \mathbf{w}_? \frac{2}{2} + \frac{2}{k} \\ 2\text{Risk}_{k, \text{train}} + 2 \frac{2}{k} &= 2\text{Risk}_{k, \text{train}} + \frac{2}{k} ; \end{aligned}$$

where we denote $2\text{Risk}_{k, \text{train}} = E_k [\mathbf{w}_{\text{ridge}}^k(D_{\text{train}}) - \mathbf{w}_? \frac{2}{2} + \frac{2}{k}]$, and we also note that $\text{Risk}_{k, \text{train}} \leq 1 + \frac{2}{k}$ by definition. By Lemma J.1, we have

$$\text{Risk}_{k, \text{train}} = \text{BayesRisk}_k + O\left(\frac{2}{k} + 1\right) \frac{N_{\text{val}}}{N} :$$

We next deal with the term $"_{\text{val}} := \max_{l \in [K]} \mathbb{E}_{\text{val}}(\mathbf{w}_l(D_{\text{train}})) - L_{\text{val}, \mathbf{w}_?}(\mathbf{w}_l(D_{\text{train}}))$. Note that for the good event $E_{\text{train}} := E \setminus E_W \setminus E_{b, \text{train}}$ of $(\mathbf{w}_?; D_{\text{train}})$, we have

$$E_k [1fE_0g "_{\text{val}}] = E_k [1fE_{\text{train}}g "_{\text{val}}] = E_{\mathbf{w}_?; D_{\text{train}} \sim P_{\mathbf{w}_?}} [1fE_{\text{train}}g E_{\text{val}} ["_{\text{val}} \mid \mathbf{w}_?; D_{\text{train}}]] :$$

Thus, Lemma J.2 yields

$$E_k [1fE_0g "_{\text{val}}] = O\left(B_w^2 \frac{\log K}{N_{\text{val}}} + \frac{\log K}{N_{\text{val}}}\right) :$$

Therefore, we can conclude that

$$E_k [1fEg(\mathbf{y}_{N+1} - y_{N+1})^2] \leq 2\text{BayesRisk}_k + O\left(\frac{2}{k} + 1\right) \frac{N_{\text{val}}}{N} + B_w^2 \frac{\log K}{N_{\text{val}}} + \frac{B_w^2 \log K}{N_{\text{val}}} :$$

Therefore, we can choose (N_{val}) so that $N_{\text{val}} = N/2$ as

$$N_{\text{val}} = \max\left(\frac{B_w^2}{2} N, \frac{2}{\log^{1=3}(K)} \log K\right) ; \quad " = \frac{\max}{N} :$$

It is worth noting that such choice of N_{val} is feasible as long as $N \gtrsim \frac{B_w^4}{4_{\text{max}}} \log K$. Under such choice, we obtain

$$\frac{1}{2} E_k [1fEg(\mathbf{y}_{N+1} - y_{N+1})^2] \leq \text{BayesRisk}_k + O\left(\frac{4=3}{\max} B_w^{2=3} \frac{\log K}{N}\right) :$$

Part II. Similar to the proof of Corollary 6, we have

$$E \frac{1}{N} \sum_{i=1}^N (y_{N+1} - \hat{y}_{N+1})^2 = O\left(\frac{B_y^2}{N^5}\right) = O\left(\frac{2}{N^4}\right);$$

Conclusion. Combining the both cases, we obtain

$$\begin{aligned} E_k \frac{1}{2} (y_{N+1} - \hat{y}_{N+1})^2 &= \text{BayesRisk}_k + O\left(\frac{4-3}{\max} B_w^{2=3} \frac{\log K}{N}\right)^{1=3} \\ &= \text{BayesRisk}_k + O\left(\frac{2}{\max} \frac{\log K}{N}\right)^{1=3} + \frac{4-3}{\max} \frac{\log^{2=3}(N) \log^{1=3}(K)}{d^{2=3} N^{1=3}} \\ &= \text{BayesRisk}_k + O\left(\frac{2}{\max} \frac{\log K}{N}\right)^{1=3}; \end{aligned}$$

where we plug in our choice of B_y . The bounds on $M^{(\cdot)}; D^{(\cdot)}$ and $\int \int$ follows immediately from Theorem I.2. This completes the proof. \square

J.2 Derivation of the exact Bayes predictor

Let $(D; \mathbf{x}_{N+1}; y_{N+1})$ be $(N + 1)$ observations from the data generating model considered in Section 4.1.1. On observing $(D; \mathbf{x}_{N+1})$, the Bayes predictor of y_{N+1} is given by its posterior mean:

$$E[y_{N+1} | D; \mathbf{x}_{N+1}] = E[\mathbf{w}^\top \mathbf{x}_{N+1} | D; \mathbf{x}_{N+1}] = \mathbf{x}_{N+1}^\top E[\mathbf{w} | D];$$

It thus remains to derive $E[\mathbf{w} | D]$. Recall that our data generating model is given by k , By Bayes' rule, we have

$$E[\mathbf{w} | D] = \sum_{k' \in [K]} P(k = k' | D) E[\mathbf{w} | D; k = k']; \quad (46)$$

On $k = k'$, the data is generated from the noisy linear model $\mathbf{w} \sim N(\mathbf{0}; \mathbf{I}_d=d)$, and $\mathbf{y} = \mathbf{X}\mathbf{w} + \mathbf{u}$ where $u_i \stackrel{\text{iid}}{\sim} N(0; \frac{2}{k'})$. It is a standard result that $E[\mathbf{w} | D; k = k']$ is given by the ridge estimator

$$\begin{aligned} E[\mathbf{w} | D; k = k'] &= \frac{\mathbf{X}^\top \mathbf{X} + d \frac{2}{k'} \mathbf{I}_d^{-1}}{\left\{ \frac{\mathbf{X}^\top \mathbf{X} + d \frac{2}{k'} \mathbf{I}_d^{-1}}{b_{k'}^{-1}} \right\}} \mathbf{X}^\top \mathbf{y} =: \mathbf{w}_{k'} \\ &= \frac{\mathbf{X}^\top \mathbf{X}}{N} + \frac{d \frac{2}{k'}}{N} \mathbf{I}_d^{-1} \frac{\mathbf{X}^\top \mathbf{y}}{N}. \end{aligned}$$

(Note that the sample covariance within $b_{k'}$ is not normalized by N , which is not to be confused with remaining parts within the paper.) Therefore, the posterior mean (46) is exactly a weighted combination of K ridge regression estimators, each with regularization $d \frac{2}{k'} = N$.

It remains to derive the mixing weights $P(k = k' | D)$ for all $k' \in [K]$. By Bayes' rule, we have

$$\begin{aligned} P(k = k' | D) &= \frac{P(k = k') \int_{\mathbf{w}} p(\mathbf{w}) p_{k', \mathbf{w}}(D | \mathbf{w}) d\mathbf{w}}{\int_{k'} \frac{1}{(2 \frac{2}{k'})^{d=2} (2 \frac{2}{k'})^{N=2}} \exp\left\{-\frac{d\mathbf{w}^\top \mathbf{w}}{2} - \frac{k \mathbf{X} \mathbf{w} - \mathbf{y}}{2 \frac{2}{k'}}\right\} d\mathbf{w}} \\ &= \frac{1}{(2 \frac{2}{k'})^{N=2}} \exp\left\{-\frac{1}{2} \mathbf{w}^\top \left(\frac{\mathbf{X}^\top \mathbf{X}}{k'} + d \mathbf{I}_d\right) \mathbf{w} + \frac{\mathbf{X}^\top \mathbf{y}}{2 \frac{2}{k'}}\right\} \frac{1}{(2 \frac{2}{k'})^{d=2}} \\ &= \frac{1}{(2 \frac{2}{k'})^{N=2}} \exp\left\{-\frac{1}{2 \frac{2}{k'}} (\mathbf{w} - \mathbf{w}_{k'})^\top b_{k'} (\mathbf{w} - \mathbf{w}_{k'})\right\} \frac{1}{(2 \frac{2}{k'})^{d=2}} \frac{1}{(2 \frac{2}{k'})^{N=2}} \mathbf{y}^\top \mathbf{X} b_{k'}^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \frac{\det(b_{k'} = \frac{2}{k'})^{-1=2}}{N \frac{2}{k'}} \exp\left\{-\frac{1}{2 \frac{2}{k'}} \mathbf{y}^\top \mathbf{X} b_{k'}^{-1} \mathbf{X}^\top \mathbf{y}\right\} \end{aligned}$$

$$\frac{1}{k^0} \frac{1}{N^{-d} \det(\mathbf{X}^\top \mathbf{X} + d \frac{2}{k^0} \mathbf{I}_d)^{1=2}} \exp \left(-\frac{1}{2} \frac{1}{k^0} \mathbf{y}^\top \mathbf{X} \mathbf{w}_{k^0} \right) :$$

Note that such mixing weights involve the determinant of the matrix $\mathbf{b}_{k^0} = \mathbf{X}^\top \mathbf{X} + d \frac{2}{k^0} \mathbf{I}_d$, which depends on the data \mathbf{X} in a non-trivial fashion; Any transformer has to approximate these weights if their mechanism is to directly approximate the exact Bayesian predictor (46).

J.3 Useful lemmas

Lemma J.1. For $2\text{Risk}_{k;\text{train}} = \mathbb{E}_k \mathbf{w}_{\text{ridge}}^k(D_{\text{train}}) \mathbf{w}_{\frac{2}{k}} + \frac{2}{k}$, there exists universal constant C such that

$$\text{Risk}_{k;\text{train}} \leq \text{BayesRisk}_k + C \left(\frac{2}{k} + 1 \right) \frac{N_{\text{val}}}{N} :$$

Proof. Recall that under \mathbb{P}_k , we have

$$\mathbf{w}_{\frac{2}{k}} \sim \mathcal{N}(0; \mathbf{I}_{d=d}); \quad y_i = \mathbf{h}_{\mathbf{x}_i} \mathbf{w}_{\frac{2}{k}} + \epsilon_i; \quad \epsilon_i \sim \mathcal{N}(0; 2) :$$

We denote $\mathbf{y}_t = [y_i]_{i \in \mathcal{I}_{\text{train}}}$, then by definition $\mathbf{w}_{\text{ridge}}^k(D_{\text{train}}) = (\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} + d \frac{2}{k})^{-1} \mathbf{X}_{\text{train}} \mathbf{y}_t$ (with $k = d \frac{2}{k} = N_{\text{train}}$). Thus, a simple calculation yields

$$2\text{Risk}_{k;\text{train}} = \mathbb{E}_k \mathbf{w}_{\text{ridge}}^k(D_{\text{train}}) \mathbf{w}_{\frac{2}{k}} + \frac{2}{k} = \frac{2}{k} \mathbb{E} \text{tr} \left(\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} + d \frac{2}{k} \right)^{-1} + \frac{2}{k} ;$$

and analogously, $2\text{BayesRisk}_k = \frac{2}{k} \mathbb{E} \text{tr} \left(\mathbf{X}^\top \mathbf{X} + d \frac{2}{k} \mathbf{I}_d \right)^{-1} + \frac{2}{k}$. Therefore,

$$2\text{Risk}_{k;\text{train}} - 2\text{BayesRisk}_k = \frac{2}{k} \mathbb{E} \text{tr} \left(\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} + d \frac{2}{k} \mathbf{I}_d \right)^{-1} - \frac{2}{k} \mathbb{E} \text{tr} \left(\mathbf{X}^\top \mathbf{X} + d \frac{2}{k} \mathbf{I}_d \right)^{-1} \\ = \frac{2}{k} N_{\text{val}} \mathbb{E}_k \left[\min \left(\frac{d}{N_{\text{train}}} \right)^{-1} \right];$$

where in the above inequality we denote $\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} + d \frac{2}{k} \mathbf{I}_d$ and use the following fact:

$$\text{tr} \left(\mathbf{I}_d + \mathbf{X}_v^\top \mathbf{X}_v \right)^{-1} = \text{tr} \left(\mathbf{I}_d + \left(\mathbf{I}_d + \mathbf{X}_v^\top \mathbf{X}_v \right)^{-1} \right)^{-1} \\ = \text{tr} \left(\mathbf{I}_d + \mathbf{X}_v^\top \mathbf{X}_v \right)^{-1} \left(\mathbf{I}_d + \mathbf{X}_v^\top \mathbf{X}_v \right)^{-1} \\ = \mathbb{E} \left[\left(\mathbf{I}_d + \mathbf{X}_v^\top \mathbf{X}_v \right)^{-1} \left(\mathbf{I}_d + \mathbf{X}_v^\top \mathbf{X}_v \right)^{-1} \right] \\ = \frac{\text{rank} \left(\mathbf{X}_v^\top \mathbf{X}_v \right)}{\max \left(\mathbf{I}_d \right)} N_{\text{val}} \min \left(\frac{d}{N_{\text{train}}} \right)^{-1} ;$$

Case 1. We first suppose that $N_{\text{train}} \geq 16d$. Then by definition $\frac{d}{N_{\text{train}}} \geq \frac{1}{16}$, and hence

$$\frac{2}{k} N_{\text{val}} \mathbb{E}_k \left[\min \left(\frac{d}{N_{\text{train}}} \right)^{-1} \right] \leq \frac{2}{k} \frac{N_{\text{val}}}{d} \frac{16 N_{\text{val}}}{N_{\text{train}}} \leq \frac{32 N_{\text{val}}}{N} ;$$

Case 2. When $N_{\text{train}} < 9d$, then we consider the event $E_t := \mathbb{P} \left(\min \left(\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} \right) = N_{\text{train}} \right) \geq \frac{1}{16} g$. By Lemma B.2 we have $\mathbb{P}(E_t^c) \leq \exp(-N_{\text{train}} g)$. Therefore,

$$\frac{2}{k} N_{\text{val}} \mathbb{E}_k \left[\min \left(\frac{d}{N_{\text{train}}} \right)^{-1} \right] = \frac{2}{k} N_{\text{val}} \mathbb{E}_k \left[1_{E_t} \min \left(\frac{d}{N_{\text{train}}} \right)^{-1} \right] + \frac{2}{k} N_{\text{val}} \mathbb{E}_k \left[1_{E_t^c} \min \left(\frac{d}{N_{\text{train}}} \right)^{-1} \right] \\ \leq \frac{16}{N_{\text{train}}} \frac{2}{k} N_{\text{val}} \mathbb{P}(E_t) + \frac{N_{\text{val}}}{d} \mathbb{P}(E_t^c) \\ \leq \frac{32}{N} \frac{2}{k} N_{\text{val}} + \frac{N_{\text{val}}}{d} \exp(-N_{\text{train}} g) = O \left(\frac{(2/k + 1) N_{\text{val}}}{N} \right) ;$$

Combining these two cases finishes the proof. \square

Lemma J.2. Condition on the event E_{train} , we have

$$\mathbb{E}_{D_{\text{val}} \sim \mathbb{P}_k | \mathbf{w}_{\frac{2}{k}}; D_{\text{train}}} \max_{J \in [K]} \mathbb{E}_{\text{val}}(\mathbf{w}_J) - L_{\text{val}; \mathbf{w}_{\frac{2}{k}}}(\mathbf{w}_J) \leq C B_w^2 \frac{\log(2K)}{N_{\text{val}}} + \frac{\log(2K)}{N_{\text{val}}} ;$$

where we denote $\mathbf{w}_J = \mathbf{w}_J(D_{\text{train}})$.

Proof. We only need to work with a fixed pair of $(\mathbf{w}_\gamma; D_{\text{train}})$ such that E_{train} holds. Hence, in the following we only consider the randomness of D_{val} conditional on such a $(\mathbf{w}_\gamma; D_{\text{train}})$.

Recall that for any \mathbf{w} ,

$$\mathbb{E}_{\text{val}}(\mathbf{w}) = \frac{1}{2jD_{\text{val}}} \times \sum_{(\mathbf{x}_i; y_i) \in D_{\text{val}}} (\mathbf{h}\mathbf{x}_i; \mathbf{w} \mid y_i)^2;$$

and we have $E_{D_{\text{val}}}[\mathbb{E}_{\text{val}}(\mathbf{w})] = L_{\text{val}; \mathbf{w}_\gamma}(\mathbf{w})$. For each $i \geq 1$,

$$y_i - \mathbf{h}\mathbf{x}_i; \mathbf{w}_\gamma \mid y_i = \sum_{l \in [K]} \mathbf{h}\mathbf{x}_i; \mathbf{w}_l - \mathbf{w}_\gamma \mid y_i \sim \text{SG}\left(\frac{2}{k} + k\mathbf{w}_\gamma, \sum_{l \in [K]} \mathbf{w}_l k^2\right);$$

Note that under E_{train} , we have $\mathbf{w}_l \geq B_2(B_W)$ for all $l \in [K]$, and hence $\frac{2}{k} + k\mathbf{w}_\gamma \leq 5B_W^2$. We then have $(y_i - \mathbf{h}\mathbf{x}_i; \mathbf{w}_l)^2$'s are (conditional) i.i.d random variables in $\text{SE}(CB_W^4)$. Then, by Bernstein's inequality, we have

$$P_{D_{\text{val}}}[\mathbb{E}_{\text{val}}(\mathbf{w}_l) - L_{\text{val}; \mathbf{w}_\gamma}(\mathbf{w}_l) \geq t] \leq 2 \exp(-cN_{\text{val}} \min\left\{\frac{t^2}{B_W^2}, \frac{t}{B_W}\right\});$$

where c is a universal constant. Applying the union bound, we obtain

$$P_{D_{\text{val}}} \max_{l \in [K]} [\mathbb{E}_{\text{val}}(\mathbf{w}_l) - L_{\text{val}; \mathbf{w}_\gamma}(\mathbf{w}_l) \geq t] \leq K \exp(-cN_{\text{val}} \min\left\{\frac{t^2}{B_W^2}, \frac{t}{B_W}\right\});$$

Taking integration completes the proof. \square

J.4 Generalized linear models with adaptive link function selection

Suppose that $(g_k : \mathbb{R} \rightarrow \mathbb{R})_{k \in [K]}$ is a set of link functions such that g_k is non-decreasing and C^2 -smooth for each $k \in [K]$. We consider the input format we introduce in Section 4.1 with $jD_{\text{train}} = dN = 2e; jD_{\text{val}} = bN = 2c$.

Theorem J.2 (GLMs with adaptive link function selection). *For any fixed set of parameters defined in Assumption B, as long as $N = O(d)$, there exists a transformer with $L = O(\log(N))$ layers, input dimension $D = O(dK)$ and $\max_{\cdot \in [L]} M(\cdot) \leq \delta^{\beta} N$, such that the following holds.*

For any $k \in [K]$ and any distribution \mathbb{P} that is a generalized linear model of the link function g_k and some parameter γ , if Assumption B holds for each pair $(\mathbb{P}; g_k)$, then

$$E_{(D; \mathbf{x}_{N+1}; y_{N+1}) \sim \mathbb{P}} (\mathbb{E}_{N+1} - y_{N+1})^2 = E_{(\mathbf{x}; y) \sim \mathbb{P}} (g_k(\mathbf{h}\mathbf{x}; \gamma) - y)^2 + O\left(\frac{d}{N} + \frac{\log(K)}{N}\right);$$

or equivalently, $E_{(D; \mathbf{x}_{N+1}) \sim \mathbb{P}} [(\mathbb{E}_{N+1} - E[y_{N+1} | \mathbf{x}_{N+1}])^2] = O\left(\frac{d}{N} + \frac{\log(K)}{N}\right)$.

Proof. For each $k \in [K]$, we consider optimizing the following training loss:

$$\mathbf{w}_{\text{GLM}}^{(k)} := \arg \min_{\mathbf{w}} \mathbb{E}_{\text{train}}^{(k)}(\mathbf{w}) := \frac{1}{N_{\text{train}}} \sum_{(\mathbf{x}_i; y_i) \in D_{\text{train}}} \ell_k(\mathbf{h}\mathbf{x}_i; \mathbf{w} \mid y_i);$$

where $\ell_k(t; y) := y^2 + \int_0^t g_k(s) ds$ is the convex (integral) loss associated with g_k (as in Section 3.1).

Also, for each predictor $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we consider the squared validation loss \mathbb{E}_{val} :

$$\mathbb{E}_{\text{val}}(f) := \frac{1}{2N_{\text{val}}} \sum_{(\mathbf{x}_i; y_i) \in D_{\text{val}}} (f(\mathbf{x}_i) - y_i)^2;$$

Fix a large universal constant C_0 . Let us set

$$\begin{aligned} B_x &= C_0 K_x \sqrt{\frac{d \log(N)}{N}}; & B_y &= C_0 K_y \sqrt{\frac{d \log(N)}{N}}; \\ B_x &= C_0 K_x \sqrt{8}; & B_y &= 8L_g K_y; \end{aligned}$$

Then, we define good events similarly to the proof of Corollary 6 (Appendix F.4):

$$\begin{aligned}
 E_W &= \bigcap_{k \geq [K]; \delta w \geq B_2(B_w); \min(r^{-2} \mathbb{E}_{\text{train}}^{(k)}(\mathbf{w})) \leq \max(r^{-2} \mathbb{E}_{\text{train}}^{(k)}(\mathbf{w})) \leq \epsilon; \\
 E_r &= \bigcap_{k \geq [K]; \mathbf{w}_{\text{GLM}}^{(k)} \in B_{w=2}}; \\
 E_{b,\text{train}} &= \bigcap_{f \in \mathcal{F}; \mathbf{x}_i; y_i \in D_{\text{train}}; k \geq k_2} B_{\mathbf{x};j} y_{ij} B_{y,g}; \\
 E_{b,\text{val}} &= \bigcap_{f \in \mathcal{F}; \mathbf{x}_i; y_i \in D_{\text{val}}; k \geq k_2} B_{\mathbf{x};j} y_{ij} B_{y,g}; \\
 E_{b;N+1} &= \bigcap_{f \in \mathcal{F}; \mathbf{x}_{N+1}; y_{N+1}} B_{\mathbf{x};j} y_{N+1,j} B_{y,g};
 \end{aligned}$$

Similar to the proof of Theorem G.2 (Appendix G.2), we know the good event $E := E_W \setminus E_r \setminus E_{b,\text{train}} \setminus E_{b,\text{test}} \setminus E_{b;N+1}$ holds with high probability: $\mathbb{P}(E^c) = O(N^{-10})$.

Similar to the proof of Theorem I.2, we can show that there exists a transformer with prediction $\hat{y}_{N+1} = \text{ready}(\text{TF}(\mathbf{H}))$ (clipped by B_y), such that (for any P) the following holds under E :

- (a) For each $k \geq [K]$, $f_k = A_k(D_{\text{train}})$ is a predictor such that $f_k(\mathbf{x}_i) = g_k(\mathbf{h}(\mathbf{x}_i; \mathbf{w}_{\text{GLM}}^{(k)}))$ for all $i \geq [N+1]$ (where “ is chosen as in Appendix G.2).
- (b) $\hat{y}_{N+1} = \text{clip}_{B_y}(\hat{y}(\mathbf{x}_{N+1}))$, where $\hat{y} = A_{\text{TF}}(D)$ is an aggregated predictor given by $\hat{y} = \sum_{k \in [K]} \lambda_k f_k$, such that (λ_k) is a distribution supported on $k \geq [K]$ such that $\mathbb{E}_{\text{val}}(f_k) \leq \min_{k \in [K]} \mathbb{E}_{\text{val}}(f_k) + \epsilon$.

Similar to the proof of Theorem G.2, for $E_0 := E_W \setminus E_r \setminus E_{b,\text{train}} \setminus E_{b,\text{test}}$, we have

$$\mathbb{E}_{(D; \mathbf{x}_{N+1}; y_{N+1}) \sim P} (\hat{y}_{N+1} - y_{N+1})^2 \leq \mathbb{E}_{D \sim P} [1fE_0 g L_{\text{val}}(\hat{y})] + O\left(\frac{B_y^2}{N^5}\right);$$

where we denote $L_{\text{val}}(f) := \mathbb{E}_{(\mathbf{x}; y) \sim P} [1f k \mathbf{x} k_2 B_{\mathbf{x};j} g(f(\mathbf{x}) - y)^2]$ for each predictor f . By the definition of \hat{y} , we then have (under E_0)

$$L_{\text{val}}(\hat{y}) \leq L_{\text{val}}(f_{k^*}) + \max_j L_{\text{val}}(f_j) - L_{\text{val}}(f_j) + \epsilon;$$

For the first term, repeating the argument in the proof of Theorem G.2 directly yields that for $E_{\text{train}} := E_W \setminus E_r \setminus E_{b,\text{train}}$,

$$\mathbb{E}_{D_{\text{train}} \sim P} [1fE_{\text{train}} g L_{\text{val}}(f_{k^*})] \leq \mathbb{E}_{(\mathbf{x}; y) \sim P} (g_{k^*}(\mathbf{h}(\mathbf{x}; i)) - y)^2 + O(d=N_{\text{train}});$$

For the second term, similar to Lemma J.2, we can show that conditional on D_{train} such that E_{train} holds, it holds

$$\mathbb{E}_{D_{\text{val}} \sim P | D_{\text{train}}} [1fE_0 g \max_j L_{\text{val}}(f_j) - L_{\text{val}}(f_j)] \leq O(K^2) \left(\frac{\log K}{N_{\text{val}}} + \frac{\log K}{N_{\text{val}}} \right);$$

Combining these inequalities and suitably choosing ϵ complete the proof. \square

K Analysis of pretraining

Thus far, we have established the existence of transformers for performing various ICL tasks with good in-context statistical performance. We now analyze the sample complexity of pretraining these transformers from a finite number of training ICL instances.

K.1 Generalization guarantee for pretraining

Setup At pretraining time, each training ICL instance has form $\mathbf{Z} := (\mathbf{H}; y_{N+1})$, where $\mathbf{H} := \mathbf{H}(D; \mathbf{x}_{N+1}) \in \mathbb{R}^{D \times (N+1)}$ denote the input sequence formatted as in (3). We consider the square loss between the in-context prediction and the ground truth label:

$$\ell_{\text{icl}}(\cdot; \mathbf{Z}) := \frac{1}{2} \left(y_{N+1} - \underbrace{\text{clip}_{B_y} \text{ready}(\text{TF}^R(\mathbf{H}))}_{\text{ready}} \right)^2;$$

Above, $\text{clip}_{B_y}(t) := \max\{f \min\{t; B_y g\}; B_y g\}$ is the standard clipping operator onto $[B_y; B_y]$, and TF^R the transformer architecture as in Definition 3 with clipping operators after each layer: let $\mathbf{H}^{(0)} = \text{clip}_R(\mathbf{H})$,

$$\mathbf{H}^{(\ell)} = \text{clip}_R \left(\text{MLP}_{\text{mlp}}^{(\ell)} \left(\text{Attn}_{\text{attn}}^{(\ell)} \left(\mathbf{H}^{(\ell-1)} \right) \right) \right) \quad \text{for all } \ell \geq [L]; \quad \text{clip}_R(\mathbf{H}) := [\text{Proj}_{\|\mathbf{h}\|_2 \leq R}(\mathbf{h}_i)]_i$$

The clipping operator is used to control the Lipschitz constant of TF with respect to \mathbf{h} , and we typically choose a sufficiently large clipping radius R so that it does not modify the behavior of the transformer on any input sequence of our concern.

We draw ICL instances $\mathbf{Z} := (\mathbf{H}; y_{N+1}) = (D; (\mathbf{x}_{N+1}; y_{N+1}))$ from a (meta-)distribution denoted as \mathcal{P} , which first sample an in-context data distribution $\mathcal{P}^{\otimes(N+1)}$, then sample iid examples $(\mathbf{x}_i; y_i)_{i=1}^{N+1} \stackrel{\text{iid}}{\sim} \mathcal{P}^{\otimes(N+1)}$ and form $D = f(\mathbf{x}_i; y_i)_{g_{i \in [N]}}$. Our pretraining loss is the average ICL loss on n pretraining instances $\mathbf{Z}^{(1:n)} \stackrel{\text{iid}}{\sim} \mathcal{P}$, and we consider the corresponding test ICL loss on a new test instance:

$$\hat{\mathcal{L}}_{\text{icl}}(\cdot) := \frac{1}{n} \sum_{j=1}^n \mathcal{L}_{\text{icl}}(\cdot; \mathbf{Z}^j); \quad \mathcal{L}_{\text{icl}}(\cdot) := \mathbb{E}_{\mathbf{P} \sim \mathcal{P}; \mathbf{Z} \sim \mathcal{P}^{(N+1)}} [\mathcal{L}_{\text{icl}}(\cdot; \mathbf{Z})]$$

Our pretraining algorithm is to solve a standard constrained empirical risk minimization (ERM) problem over transformers with L layers, M heads, and norm bound B (recall the definition of the $\|\cdot\|$ norm in (2)):

$$\begin{aligned} \mathbf{b} &:= \arg \min_{\in \mathcal{L}; M; D^0; B} \hat{\mathcal{L}}_{\text{icl}}(\cdot); \\ \mathcal{L}; M; D^0; B &:= \left(\begin{matrix} (1:L) \\ \text{attn} \end{matrix}; \begin{matrix} (1:L) \\ \text{mlp} \end{matrix} \right); \max_{\in [L]} M^{(\cdot)}; \max_{\in [L]} D^{(\cdot)}; D'; \|\cdot\|; B \end{aligned} \quad (\text{TF-ERM})$$

Generalization guarantee By standard uniform concentration analysis via chaining arguments (Proposition B.4; see also [87, Chapter 5] for similar arguments), we have the following excess loss guarantee for (TF-ERM). The proof can be found in Appendix L.2.

Theorem K.1 (Generalization for pretraining). *With probability at least $1 - \epsilon$ (over the pretraining instances $f\mathbf{Z}^j g_{j \in [n]}$), the solution \mathbf{b} to (TF-ERM) satisfies*

$$\mathcal{L}_{\text{icl}}(\mathbf{b}) \leq \inf_{\in \mathcal{L}; M; D^0; B} \mathcal{L}_{\text{icl}}(\cdot) + O \left(B_y^2 \frac{\log(1/\epsilon)}{n} \right);$$

where $\log = \log(2 + \max\{fB; R; B_y g\})$ is a log factor.

K.2 Examples of pretraining for in-context regression problems

In Theorem K.1, the comparator $\inf_{\in \mathcal{L}; M; D^0; B} \mathcal{L}_{\text{icl}}(\cdot)$ is simply the smallest expected ICL loss for ICL instances drawn from \mathcal{P} , among all transformers within the norm ball $\mathcal{L}; M; D^0; B$. Using our constructions in Section 3 & 4, we show that this comparator loss is small on various (meta-)distribution \mathcal{P} 's, by which we obtain end-to-end guarantees for pretraining transformers with small ICL loss at test time. Here we showcase this argument on several representative regression problems.

Linear regression For any in-context data distribution \mathcal{P} , let $\mathbf{w}_P^2 := \mathbb{E}_P[\mathbf{x}\mathbf{x}^\top]^{-1} \mathbb{E}_P[\mathbf{x}y]$ denote the best linear predictor for \mathcal{P} . We show that with mild choices of $\mathcal{L}; M; B$, the learned transformer can perform in-context linear regression with near-optimal statistical power, in that on the sampled \mathcal{P} and ICL instance $f(\mathbf{x}_i; y_i)_{g_{i \in [N+1]}} \stackrel{\text{iid}}{\sim} \mathcal{P}$, it competes with the best linear predictor \mathbf{w}_P^2 for this particular \mathcal{P} . The proof follows directly by combining Corollary 5 with Theorem K.1, and can be found in Appendix L.3.

Theorem K.2 (Pretraining transformers for in-context linear regression). *Suppose \mathcal{P} is almost surely well-posed for in-context linear regression (Assumption A) with the canonical parameters. Then, for $N \in \Theta(d)$, with probability at least $1 - \epsilon$ (over the training instances $\mathbf{Z}^{(1:n)}$), the solution*

\mathbb{P} of (TF-ERM) with $L = O(\log(N))$ layers, $M = 3$ heads, $D' = 0$ (attention-only), and $B = O(\frac{1}{d})$ achieves small excess ICL risk over $\mathbb{P}_P^?$:

$$L_{\text{icl}}(\mathbb{b}) = \mathbb{E}_{\mathbb{P} \sim \mathbb{P}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}} \frac{1}{2} (y - h_{\mathbb{P}_P^?}(\mathbf{x}))^2 \in \Theta \left(\frac{\sqrt{2d^2 + \log(1/n)}}{n} + \frac{d^2}{N} \right);$$

where $\Theta(\cdot)$ only hides polylogarithmic factors in $d; N; 1/n$.

To our best knowledge, Theorem K.2 offers the first end-to-end result for pretraining a transformer to perform in-context linear regression with explicit excess loss bounds. The $\Theta(\frac{1}{\sqrt{2d^2+n}})$ term originates from the generalization of pretraining (Theorem K.1), where as the $\Theta(d^2=N)$ term agrees with the standard fast rate for the excess loss of linear regression [38]. Further, as long as $n \in \Theta(2N^2)$, the excess risk achieves the optimal rate $\Theta(d^2=N)$ (up to log factors).

Additional examples By similar arguments as in the proof of Theorem K.2, we can directly turn most of our other expressivity results into results on the pretrained transformers. Here we present three such additional examples (proofs in Appendix L.4-L.6). The first example is for the sparse linear regression problem considered in Theorem 8.

Theorem K.3 (Pretraining transformers for in-context sparse linear regression). *Suppose each \mathbb{P} is almost surely an instance of the sparse linear model specified in Theorem 8 with parameters $B_W^?$ and s . Suppose $N \in \Theta(s \log((d-N)))$ and let $\mathbb{P} := B_W^? =$.*

Then with probability at least $1 - \frac{1}{n}$ (over the training instances $\mathbf{Z}^{(1:n)}$), the solution \mathbb{b} of (TF-ERM) with $L = \Theta(2(1+dN))$ layers, $M = 2$ heads, $D' = 2d$, and $B = \Theta(\text{poly}(d; B_W^?))$ achieves small excess ICL risk:

$$L_{\text{icl}}(\mathbb{b}) \leq \frac{1}{2} \in \Theta \left(\frac{\sqrt{4d^2(1+dN)^2 + \log(1/n)}}{n} + \frac{2s \log d}{N} \right);$$

where $\Theta(\cdot)$ only hides polylogarithmic factors in $d; N; 1/n$.

Our next example is for the problem of noisy linear regression with mixed noise levels considered in Theorem 12 and Theorem J.1. There, the constructed transformer uses the post-ICL validation mechanism to perform ridge regression with an adaptive regularization strength depending on the particular input sequence.

Theorem K.4 (Pretraining transformers for in-context noisy linear regression with algorithm selection). *Suppose \mathbb{P} is the data generating model (noisy linear model with mixed noise levels) considered in Theorem J.1, with $\max_{\mathbf{x}} \|\mathbf{x}\| \leq O(1)$. Let $N \geq d=10$.*

Then, with probability at least $1 - \frac{1}{n}$ (over the training instances $\mathbf{Z}^{(1:n)}$), the solution \mathbb{b} of (TF-ERM) with input dimension $D = O(dK)$, $L = O(\frac{1}{\min} \log(N/\min))$ layers, $M = O(K)$ heads, $D' = O(K^2)$, and $B = O(\text{poly}(K; \frac{1}{\min}; d; N))$ achieves small excess ICL risk:

$$L_{\text{icl}}(\mathbb{b}) \leq \text{BayesRisk} \in \Theta \left(\frac{\sqrt{\frac{-4}{\min} K^3 d^2 + \log(1/n)}}{n} + \frac{2 \max_{\mathbf{x}} \log K}{N} \right)^{1-3} A;$$

where $\Theta(\cdot)$ only hides polylogarithmic factors in $d; N; K; 1/\min$.

Our final example is for in-context logistic regression. For simplicity we consider the realizable case.

Theorem K.5 (Pretraining transformers for in-context logistic regression; square loss guarantee). *Suppose for \mathbb{P} , \mathbb{P} is almost surely a realizable logistic model (i.e. $\mathbb{P} = \mathbb{P}^{\log}$ with $k_2 \leq B_W^?$ as in Corollary G.1). Suppose that $B_W^? = O(1)$ and $N \geq O(d)$.*

Then, with probability at least $1 - \frac{1}{n}$ (over the training instances $\mathbf{Z}^{(1:n)}$), the solution \mathbb{b} of (TF-ERM) with $L = O(\log(N))$ layers, $M = \Theta(d^3 N)$ heads, $D' = 0$, and $B = O(\text{poly}(d; N))$ achieves small excess ICL risk:

$$L_{\text{icl}}(\mathbb{b}) = \mathbb{E}_{\mathbb{P}^{\log} \sim \mathbb{P}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^{\log}} \frac{1}{2} (y - \log(h(\mathbf{x}; \cdot)))^2 \in \Theta \left(\frac{\sqrt{d^6 N + \log(1/n)}}{n} + \frac{d}{N} \right);$$

where $\Theta(\cdot)$ only hides polylogarithmic factors in $d; N$.

Remark on generality of transformer All results above are established by the expressivity results in Section 3 & 4 for transformers to implement various ICL procedures (such as least squares, Lasso, GLM, and ridge regression with in-context algorithm selection), combined with the generalization bound (Theorem K.1). However, the transformer itself was not specified to encode any actual structure about the problem at hand in any result above, other than having sufficiently large number of layers, number of heads, and weight norms, which illustrates the flexibility of the transformer architecture.

L Proofs for Section K

L.1 Lipschitzness of transformers

For any $p \geq [1; 1]$, let $\|\mathbf{H}\|_{2;p} := \left(\sum_{i=1}^N \|\mathbf{h}_i\|_2^p \right)^{1/p}$ denote the column-wise $(2;p)$ -norm of \mathbf{H} . For any radius $R > 0$, we denote $H_R := \{\mathbf{H} : \|\mathbf{H}\|_{2;\infty} \leq R\}$ be the ball of radius R under norm $\|\cdot\|_{2;\infty}$.

Lemma L.1. For a single MLP layer $\text{mlp} = (\mathbf{W}_1; \mathbf{W}_2)$, we introduce its norm (as in (2))

$$\|\text{mlp}\| := \|\mathbf{W}_1\|_{\text{op}} + \|\mathbf{W}_2\|_{\text{op}}.$$

For any fixed hidden dimension D , we consider

$$\|\text{mlp}; B\| := \|\text{mlp}\| : \|\text{mlp}\| \leq B.$$

Then for $\mathbf{H} \in H_R$, $\|\text{mlp}\| \leq B$, the function $(\|\text{mlp}\|; \mathbf{H}) \mapsto \text{MLP}_{\text{mlp}}(\mathbf{H})$ is (BR) -Lipschitz w.r.t. $\|\text{mlp}\|$ and $(1 + B^2)$ -Lipschitz w.r.t. \mathbf{H} .

Proof. Recall that by our definition, for the parameter $\text{mlp} = (\mathbf{W}_1; \mathbf{W}_2) \in \|\text{mlp}; B\|$ and the input $\mathbf{H} = [\mathbf{h}_i] \in \mathbb{R}^{D \times N}$, the output $\text{MLP}_{\text{mlp}}(\mathbf{H}) = \mathbf{H} + \mathbf{W}_2 (\mathbf{W}_1 \mathbf{H}) = [\mathbf{h}_i + \mathbf{W}_2 (\mathbf{W}_1 \mathbf{h}_i)]_i$. Therefore, for $\text{mlp}' = (\mathbf{W}'_1; \mathbf{W}'_2) \in \|\text{mlp}; B\|$, we have

$$\begin{aligned} & \|\text{MLP}_{\text{mlp}}(\mathbf{H}) - \text{MLP}_{\text{mlp}'}(\mathbf{H})\|_{2;\infty} \\ &= \max_i \|\mathbf{W}_2 (\mathbf{W}_1 \mathbf{h}_i) - \mathbf{W}'_2 (\mathbf{W}'_1 \mathbf{h}_i)\|_2 \\ &= \max_i \|\mathbf{W}_2 (\mathbf{W}_1 \mathbf{h}_i) - \mathbf{W}'_2 (\mathbf{W}_1 \mathbf{h}_i) + \mathbf{W}'_2 (\mathbf{W}_1 \mathbf{h}_i) - \mathbf{W}'_2 (\mathbf{W}'_1 \mathbf{h}_i)\|_2 \\ &\leq \max_i \|\mathbf{W}_2 - \mathbf{W}'_2\|_{\text{op}} \|\mathbf{W}_1 \mathbf{h}_i\|_2 + \|\mathbf{W}'_2\|_{\text{op}} \|\mathbf{W}_1 \mathbf{h}_i - \mathbf{W}'_1 \mathbf{h}_i\|_2 \\ &\leq \max_i \|\mathbf{W}_2 - \mathbf{W}'_2\|_{\text{op}} \|\mathbf{W}_1 \mathbf{h}_i\|_2 + \|\mathbf{W}'_2\|_{\text{op}} \|\mathbf{W}_1\|_{\text{op}} \|\mathbf{h}_i - \mathbf{h}'_i\|_2 \\ &\leq BR \|\mathbf{W}_2 - \mathbf{W}'_2\|_{\text{op}} + BR \|\mathbf{W}_1 - \mathbf{W}'_1\|_{\text{op}}; \end{aligned}$$

where the second inequality follows from the 1-Lipschitzness of $\|\cdot\|_+ = [\cdot]_+$. Similarly, for $\mathbf{H}' = [\mathbf{h}'_i] \in \mathbb{R}^{D \times N}$,

$$\begin{aligned} \|\text{MLP}_{\text{mlp}}(\mathbf{H}) - \text{MLP}_{\text{mlp}}(\mathbf{H}')\|_{2;\infty} &= \max_i \|\mathbf{h}_i + \mathbf{W}_2 (\mathbf{W}_1 \mathbf{h}_i) - \mathbf{h}'_i - \mathbf{W}_2 (\mathbf{W}_1 \mathbf{h}'_i)\|_2 \\ &\leq \|\mathbf{h}_i - \mathbf{h}'_i\|_2 + \|\mathbf{W}_2 (\mathbf{W}_1 \mathbf{h}_i) - \mathbf{W}_2 (\mathbf{W}_1 \mathbf{h}'_i)\|_2 \\ &\leq \|\mathbf{h}_i - \mathbf{h}'_i\|_2 + \|\mathbf{W}_2\|_{\text{op}} \|\mathbf{W}_1 \mathbf{h}_i - \mathbf{W}_1 \mathbf{h}'_i\|_2 \\ &\leq \|\mathbf{h}_i - \mathbf{h}'_i\|_2 + B^2 \|\mathbf{h}_i - \mathbf{h}'_i\|_2; \end{aligned}$$

□

Lemma L.2. For a single attention layer $\text{attn} = f(\mathbf{V}_m; \mathbf{Q}_m; \mathbf{K}_m) g_{m \in [M]} \in \mathbb{R}^{D \times D}$, we introduce its norm (as in (2))

$$\|\text{attn}\| := \max_{m \in [M]} \|\mathbf{Q}_m\|_{\text{op}} \|\mathbf{K}_m\|_{\text{op}} + \sum_{m=1}^M \|\mathbf{V}_m\|_{\text{op}}.$$

For any fixed dimension D , we consider

$$\|\text{attn}; B\| := \|\text{attn}\| : \|\text{attn}\| \leq B.$$

Then for $\mathbf{H} \in H_R$, $\|\text{attn}\| \leq B$, the function $(\|\text{attn}\|; \mathbf{H}) \mapsto \text{Attn}_{\text{attn}}(\mathbf{H})$ is $(B^2 R^3)$ -Lipschitz w.r.t. $\|\text{attn}\|$ and $(1 + B^3 R^2)$ -Lipschitz w.r.t. \mathbf{H} .

$$\begin{aligned}
& \sum_{m=1}^M \frac{1}{N} \sum_{j=1}^N \left(\mathbf{K}_m \mathbf{V}_m \mathbf{K}_{\text{op}} \left(\mathbf{h}_{\mathbf{Q}_m} \mathbf{h}_i; \mathbf{K}_m \mathbf{h}_j \right) \right) \mathbf{h}_j - \mathbf{h}'_j \\
& \quad + \left(\mathbf{h}_{\mathbf{Q}_m} \mathbf{h}_i; \mathbf{K}_m \mathbf{h}_j \right) \left(\mathbf{h}_{\mathbf{Q}_m} \mathbf{h}'_j; \mathbf{K}_m \mathbf{h}_j \right) \mathbf{h}'_j \\
& \quad + \left(\mathbf{h}_{\mathbf{Q}_m} \mathbf{h}'_j; \mathbf{K}_m \mathbf{h}_j \right) \mathbf{Q}_m \mathbf{h}'_j; \mathbf{K}_m \mathbf{h}'_j \mathbf{h}'_j \\
& \sum_{m=1}^M \frac{1}{N} \sum_{j=1}^N \mathbf{K}_m \mathbf{V}_m \mathbf{K}_{\text{op}} \left(3 \mathbf{K}_m \mathbf{Q}_m \mathbf{K}_{\text{op}} \mathbf{K}_m \mathbf{K}_{\text{op}} R^2 \right) \mathbf{h}_j - \mathbf{h}'_j \\
& R^2 k_{\mathbf{H}} \mathbf{H}' k_{2;\infty} \left(3 \max_{m \in [M]} \mathbf{K}_m \mathbf{Q}_m \mathbf{K}_{\text{op}} \mathbf{K}_m \mathbf{K}_{\text{op}} \right) \sum_{m=1}^M \mathbf{K}_m \mathbf{V}_m \mathbf{K}_{\text{op}} \\
& B^3 R^2 k_{\mathbf{H}} \mathbf{H}' k_{2;\infty};
\end{aligned}$$

where the last inequality uses $\|\cdot\|_{\text{attn}}$ and the AM-GM inequality. This completes the proof of the Lipschitzness w.r.t. \mathbf{H} . \square

Corollary L.1. For a fixed number of heads M and hidden dimension D' , we consider

$$\text{TF}_{1;B} = (\cdot; \text{attn}; \text{mlp}) : M \text{ heads, hidden dimension } D', \|\cdot\|_{\text{attn}} \|\cdot\|_{\text{mlp}} \leq B :$$

Then for the function TF^R given by

$$\text{TF}^R : (\cdot; \mathbf{H}) \mapsto \text{clip}_R \text{MLP}_{\text{mlp}}(\text{Attn}_{\text{attn}}(\mathbf{H})) ; \|\cdot\|_{\text{attn}} \|\cdot\|_{\text{mlp}} \leq B; \mathbf{H} \in H_R$$

TF^R is B -Lipschitz w.r.t. $\|\cdot\|_{\text{attn}}$ and L_H -Lipschitz w.r.t. \mathbf{H} , where $B := BR(1 + BR^2 + B^3R^2)$ and $B_H := (1 + B^2)(1 + B^2R^3)$.

Proof. For any $(\cdot; \text{attn}; \text{mlp}), \mathbf{H} \in H_R$, and $(\cdot'; \text{attn}; \text{mlp})$, we have

$$\begin{aligned}
k_{\text{TF}}(\mathbf{H}) - \text{TF}(\mathbf{H}') k_{2;\infty} & \leq \text{MLP}_{\text{mlp}}(\text{Attn}_{\text{attn}}(\mathbf{H})) - \text{MLP}_{\text{mlp}}(\text{Attn}_{\text{attn}}(\mathbf{H}')) k_{2;\infty} \\
& \quad + \text{MLP}_{\text{mlp}}(\text{Attn}_{\text{attn}}(\mathbf{H})) - \text{MLP}_{\text{mlp}}(\text{Attn}_{\text{attn}}(\mathbf{H}')) k_{2;\infty} \\
& \quad (1 + B^2) \|\text{Attn}_{\text{attn}}(\mathbf{H}) - \text{Attn}_{\text{attn}}(\mathbf{H}')\|_{2;\infty} + B\bar{R} \|\text{mlp}\|_{\text{mlp}} \\
& \quad (1 + B^2) B^2 R^3 \|\text{attn}\|_{\text{attn}} \|\text{attn}\|_{\text{attn}} + B\bar{R} \|\text{mlp}\|_{\text{mlp}} \\
& \quad B \|\text{attn}\|_{\text{attn}} \|\text{mlp}\|_{\text{mlp}};
\end{aligned}$$

where the second inequality follows from Lemma L.2 and Lemma L.1 and the fact that $k_{\text{Attn}_{\text{attn}}(\mathbf{H})} k_{2;\infty} \bar{R} := R + B^3R^3$ for all $\mathbf{H} \in H_R$.

Furthermore, for $\mathbf{H}' \in H_R$, we have

$$\begin{aligned}
k_{\text{TF}}(\mathbf{H}) - \text{TF}(\mathbf{H}') k_{2;\infty} & \leq (1 + B^2) k_{\text{Attn}_{\text{attn}}(\mathbf{H}) - \text{Attn}_{\text{attn}}(\mathbf{H}')} k_{2;\infty} \\
& \quad (1 + B^2)(1 + B^3R^2) k_{\mathbf{H}} \mathbf{H}' k_{2;\infty};
\end{aligned}$$

which also follows from Lemma L.2 and Lemma L.1. \square

Proposition L.1 (Lipschitzness of transformers). For a fixed number of heads M and hidden dimension D' , we consider

$$\text{TF}_{L;B} = (\cdot; \text{attn}; \text{mlp}) : M^{(L)} = M; D^{(L)} = D'; \|\cdot\|_{\text{attn}} \|\cdot\|_{\text{mlp}} \leq B :$$

Then the function TF^R is $(LB_H^{L-1}B)$ -Lipschitz w.r.t. $\|\cdot\|_{\text{attn}} \|\cdot\|_{\text{mlp}} \leq B$ for any fixed \mathbf{H} .

Proof. For $(\cdot; \text{attn}; \text{mlp}) \in \text{TF}_{L;B}$, $\mathbf{e} = \mathbf{e}^{(1:L)} \in \text{TF}_{L;B}$, we have

$$\text{TF}^R(\mathbf{H}) - \text{TF}^R(\mathbf{H}') k_{2;\infty}$$

$$\begin{aligned}
& \times \sum_{l=1}^L \text{TF}_{e^{(1:l)}}^R(\mathbf{H}) - \text{TF}_{e^{(1:l)}}^R(\mathbf{H}) \quad \text{TF}_{e^{(1:l)}}^R(\mathbf{H}) - \text{TF}_{e^{(1:l)}}^R(\mathbf{H}) \quad \text{TF}_{e^{(1:l)}}^R(\mathbf{H}) - \text{TF}_{e^{(1:l)}}^R(\mathbf{H}) \quad 2; \infty \\
& \times \sum_{l=1}^L B_H^{L-l} \text{TF}_{e^{(1:l)}}^R(\mathbf{H}) - \text{TF}_{e^{(1:l)}}^R(\mathbf{H}) \quad \text{TF}_{e^{(1:l)}}^R(\mathbf{H}) - \text{TF}_{e^{(1:l)}}^R(\mathbf{H}) \quad 2; \infty \\
& \times \sum_{l=1}^L B_H^{L-l} B_H^{L-l} e^{(1:l)} - e^{(1:l)} \quad LB_H^{L-1} B_H^{L-1} e^{(1:l)} - e^{(1:l)} ;
\end{aligned}$$

where the second inequality follows from Corollary L.1, and the last inequality is because $B_H \geq 1$. \square

L.2 Proof of Theorem K.1

In this section, we prove a slightly more general result by considering the general ICL loss

$$\ell_{\text{icl}}(\cdot; \mathbf{Z}) := \ell_{\text{icl}}(\text{rad}_y(\text{TF}^R(H)); y_{N+1});$$

We assume that the loss function ℓ satisfies $\sup_j j \leq B^0$ and $\sup_{j \in \mathcal{I}} j \leq B^1$. For the special case $\ell(s; t) = \frac{1}{2}(s - t)^2$, we can take $B^0 = 4B_y$; $B^1 = 2B_y$.

We then consider

$$X := \frac{1}{n} \sum_{j=1}^n \ell_{\text{icl}}(\cdot; \mathbf{Z}^j) - \mathbb{E}_{\mathbf{Z}}[\ell_{\text{icl}}(\cdot; \mathbf{Z})];$$

where $\mathbf{Z}^{(1:n)}$ are i.i.d copies of $\mathbf{Z} \sim \mathcal{P}; \mathcal{P}$. It remains to apply Proposition B.4 to the random process fX . We verify the preconditions:

(a) By [87, Example 5.8], it holds that $\log N(\cdot; B_{\|\cdot\|}(r); \|\cdot\|) \leq L(3MD^2 + 2DD') \log(1 + 2r)$, where $B_{\|\cdot\|}(r)$ is any ball of radius r under norm $\|\cdot\|$.

(b) $j \ell_{\text{icl}}(\cdot; \mathbf{Z}) \leq B^0$ and hence B^0 -sub-Gaussian.

(c) $\ell_{\text{icl}}(\cdot; \mathbf{Z}) - \ell_{\text{icl}}(\cdot; \mathbf{e}) \leq B^1 (LB_H^{L-1} B_H^{L-1}) e$, by Proposition L.1.

Therefore, we can apply the uniform concentration result in Proposition B.4 to obtain that, with probability at least $1 - \delta$,

$$\sup_j X_j \leq CB^0 \sqrt{\frac{L(MD^2 + DD') \log(1 + 2r)}{n}};$$

where $r = \log(2 + B_H LB_H^{L-1} B_H^{L-1} B^1 = B^0) + 20L \log(2 + \max\{B; R; B^1 = B^0\}g)$. Recalling that

$$L_{\text{icl}}(b) = \inf L_{\text{icl}}(\cdot) + 2 \sup_j X_j$$

completes the proof. \square

L.3 Proof of Theorem K.2

By Corollary 5, there exists a transformer TF such that for every \mathcal{P} satisfying Assumption A with canonical parameters (and thus in expectation over \mathcal{P}) and every $N \in \Theta(d)$, it outputs prediction $y_{N+1} = \text{rad}_y(\text{TF}(\mathbf{H}))$ such that

$$L_{\text{icl}}(\cdot) = \mathbb{E}_{\mathcal{P} \sim \mathcal{P}; (\mathcal{D}; \mathbf{x}_{N+1}; y_{N+1}) \sim \mathcal{P}} \frac{1}{2} (y_{N+1} - y_{N+1})^2 = \mathbb{E}_{\mathcal{P} \sim \mathcal{P}} [L_{\mathcal{P}}(\mathbf{w}_{\mathcal{P}}^2)] + O\left(\frac{d^2}{N}\right);$$

where we recall that $L_{\mathcal{P}}(\mathbf{w}_{\mathcal{P}}^2) := \frac{1}{2} \mathbb{E}_{(\mathbf{x}; y) \sim \mathcal{P}} (y - \langle \mathbf{w}_{\mathcal{P}}^2; \mathbf{x} \rangle)^2$. By inspecting the proof, the same result holds if we change TF to the clipped version TF^R if we choose $R^2 = O(B_x^2 + B_y^2 + B_w^2 + 1) = O(d + \epsilon)$, so that on the good event $E_{\text{cov}} \setminus E_w$ considered therein, all intermediate outputs within

TF has $k_{2;\infty}$ R and thus the clipping does not modify the transformer output on $E_{\text{cov}} \setminus E_w$. Further, recall by (33) that has size bounds

$$L = O\left(\log \frac{N}{d}\right); \quad \max_{\ell \in [L]} M^{(\ell)} \leq 3; \quad \prod_{\ell \in [L]} M^{(\ell)} = O(d^{\frac{L}{2}}):$$

We can thus apply Theorem K.1 to obtain that the solution \hat{b} to (TF-ERM) with the above choice of $(L; M; B)$ and $D' = 0$ (attention-only) satisfies the following with probability at least $1 - \epsilon$:

$$L_{\text{icl}}(\hat{b}) \leq \inf_{\theta \in \mathcal{L}; M; D^0; B} L_{\text{icl}}(\theta) + O\left(\frac{L^2 MD^2 + \log(1-\epsilon)}{n}\right)$$

$$L_{\text{icl}}(\theta) + \Theta\left(\frac{L^2 MD^2 + \log(1-\epsilon)}{n}\right) \leq \Theta\left(\frac{d^{\frac{L}{2}} + \log(1-\epsilon)}{n}\right) + \frac{d^{\frac{L}{2}}}{N}:$$

Above, $\epsilon = O(\log(1 + \max_{\ell \in [L]} f_{B_y; R; B_g})) = \Theta(1)$. This finishes the proof. \square

L.4 Proof of Theorem K.3

We invoke Theorem 8 (using the construction in Theorem 7 with a different choice of L) with the following parameters:

$$L = \Theta\left(\frac{B_w^2}{d}\right)^2 = \Theta\left(\frac{1}{d}\right) \log(1 + d=N) = \Theta\left(\frac{1}{d}\right) \log(1 + d=N); \quad M = \Theta(1); \quad D' = 2d;$$

$$B_x = \Theta\left(\frac{1}{d}\right); \quad B_y = \Theta\left(\frac{1}{B_y^2} + \frac{1}{d}\right); \quad \prod_{\ell \in [L]} M^{(\ell)} = \frac{1}{B_y^2 N^2};$$

$$\prod_{\ell \in [L]} B = O\left(R + \frac{1}{N}\right)^{-1} = O\left(R + \frac{1}{\log d}\right) = \Theta(\text{poly}(d; B_w^2;));$$

where $\Theta(\cdot)$ hides polylogarithmic factors in $d; N; B_w^2$.

Then, Theorem 8 shows that there exists a transformer with L layers, $\max_{\ell \in [L]} M^{(\ell)} = M$ heads, D' hidden dimension for the MLP layers, and $\prod_{\ell \in [L]} B$ such that, on almost surely every \mathcal{P} , it returns a prediction \hat{y}_{N+1} such that, on the good event E_0 considered therein (over $D = \mathcal{P}$) which satisfies $\mathbb{P}(E_0) \geq 1 - \epsilon$,

$$\mathbb{E}_{(\mathbf{x}_{N+1}, \mathbf{y}_{N+1}) \sim \mathcal{P}} \left(\hat{y}_{N+1} - y_{N+1} \right)^2 \leq \frac{1}{2} \left[1 + O(s \log(d=N)) \right];$$

By inspecting the proof, the same result holds if we change TF to the clipped version TF^R if we choose $R^2 = O(B_x^2 + B_y^2 + (B_w^2)^2 + 1) = O(d + (B_w^2)^2 + \frac{1}{d})$, so that on the good event E_0 considered therein, all intermediate outputs within TF has $k_{2;\infty}$ R and thus the clipping does not modify the transformer output on the good event. On the bad event E_0^c , using the same argument as in the proof of Theorem 8, we have

$$\mathbb{E}_{\mathcal{D}; (\mathbf{x}_{N+1}, \mathbf{y}_{N+1}) \sim \mathcal{P}} \mathbb{1}_{E_0^c} \left(\hat{y}_{N+1} - y_{N+1} \right)^2 \leq \frac{1}{\mathbb{P}_{\mathcal{D}}(E_0)} \mathbb{E}_{\mathbf{y}_{N+1} \sim \mathcal{P}} \left(B_y^4 + y_{N+1}^4 \right)^{1/2} = \Theta\left(\frac{1}{N}\right):$$

Combining the above two bounds and further taking expectation over \mathcal{P} gives

$$L_{\text{icl}}(\hat{b}) = \mathbb{E}_{\mathcal{P} \sim \mathcal{D}; (\mathbf{x}_{N+1}, \mathbf{y}_{N+1}) \sim \mathcal{P}} \frac{1}{2} \left(\hat{y}_{N+1} - y_{N+1} \right)^2 \leq \frac{1}{2} + \Theta\left(\frac{1}{N}\right) s \log d = N:$$

We can thus apply Theorem K.1 to obtain that the solution \hat{b} to (TF-ERM) with the above choice of $(L; M; B; D')$ satisfies the following with probability at least $1 - \epsilon$:

$$L_{\text{icl}}(\hat{b}) \leq \inf_{\theta \in \mathcal{L}; M; D^0; B} L_{\text{icl}}(\theta) + O\left(\frac{L^2(MD^2 + DD') + \log(1-\epsilon)}{n}\right)$$

$$L_{\text{icl}}(\theta) + \Theta\left(\frac{L^2(MD^2 + DD') + \log(1-\epsilon)}{n}\right) \leq \Theta\left(\frac{d^{\frac{L}{2}}(1 + d=N)^2 + \log(1-\epsilon)}{n}\right) + \frac{2s \log d}{N}:$$

Above, $\epsilon = O(\log(1 + \max_{\ell \in [L]} f_{B_y; R; B_g})) = \Theta(1)$. This finishes the proof. \square

L.5 Proof of Theorem K.4

We invoke Theorem 12 and Theorem J.1, which shows that (recalling the input dimension $D = (Kd)$) there exists a transformer with the following size bounds:

$$L = O\left(\frac{1}{\min} \log(N = \min)\right); \quad \max_{\ell \in [L]} M^{(\ell)} = M = O(K); \quad \max_{\ell \in [L]} D^{(\ell)} = D' = O(K^2);$$

$$\mathbb{E} \left[\frac{1}{2} (y_{N+1} - \hat{y}_{N+1})^2 \right] = O\left(\max K d \log(N)\right);$$

such that it outputs \hat{y}_{N+1} that satisfies

$$\mathbb{E} \left[\frac{1}{2} (y_{N+1} - \hat{y}_{N+1})^2 \right] = \text{BayesRisk} + \Theta\left(\frac{2}{\min} \frac{\log K}{N} \right)^{1=3};$$

By inspecting the proof, the same result holds if we change TF to the clipped version TF^R if we choose $R^2 = O(B_x^2 + B_y^2 + (B_w^?)^2 + 1) = O(d + \frac{2}{\max})$, so that on the good event considered therein, all intermediate outputs within TF has $k k_{2,\infty} R$ and thus the clipping does not modify the transformer output on the good event. Using this clipping radius, we obtain

$$L_{\text{icl}}(\cdot) = \mathbb{E}_{P_{\sim}(\mathcal{D}; \mathbf{x}_{N+1}; y_{N+1}) \sim P} \left[\frac{1}{2} (y_{N+1} - \hat{y}_{N+1})^2 \right] = \text{BayesRisk} + \Theta\left(\frac{2}{\min} \frac{\log K}{N} \right)^{1=3};$$

We can thus apply Theorem K.1 to obtain that the solution \hat{b} to (TF-ERM) with the above choice of $(L; M; B; D')$ satisfies the following with probability at least $1 - \epsilon$:

$$L_{\text{icl}}(\hat{b}) \leq \inf_{L; M; D^0; B} L_{\text{icl}}(\cdot) + O\left(\frac{L^2(MD^2 + DD') + \log(1-\epsilon)}{n}\right)^{1=3}$$

$$L_{\text{icl}}(\cdot) + \Theta\left(\frac{L^2(MD^2 + DD') + \log(1-\epsilon)}{n}\right)^{1=3}$$

$$\leq \text{BayesRisk} + \Theta\left(\frac{\frac{-4}{\min} K^3 d^2 + \log(1-\epsilon)}{n}\right) + \frac{2}{\min} \frac{\log K}{N} \left(\frac{1}{N}\right)^{1=3} A;$$

Above, $\epsilon = O(\log(1 + \max f B_y; R; B g)) = \Theta(1)$. This finishes the proof. \square

L.6 Proof of Theorem K.5

The proof follows from similar arguments as of Theorem K.3 and Theorem K.4, where we plug in the size bounds (number of layers, heads, and weight norms) from Theorem G.2 and Corollary G.1. \square

M Experimental details and additional studies

M.1 Additional details for Section 6

Architecture and optimization We train a 12-layer encoder-only transformer, where each layer consists of an attention layer as in Definition 1 with $M = 8$ heads, hidden dimension $D = 64$, and ReLU activation (normalized by the sequence length), as well as an MLP layer as in Definition 2 hidden dimension $D' = 64$. We add Layer Normalization [3] after each attention and MLP layer to help optimization, as in standard implementations [84]. We append linear read-in layer and linear read-out layer before and after the transformer respectively, both applying a same affine transform to all tokens in the sequence and are trainable. The read-in layer maps any input vector to a D -dimensional hidden state, and the read-out layer maps a D -dimensional hidden state to a 1-dimensional scalar.

Each training sequence corresponds to a single ICL instance with N in-context training examples $f(\mathbf{x}_i; y_i) g_{i=1}^N \in \mathbb{R}^d \in \mathbb{R}$ and test input $\mathbf{x}_{N+1} \in \mathbb{R}^d$. The input to the transformer is formatted as

in (3) where each token has dimension $d + 1$ (no zero-paddings). The transformer is trained by minimizing the following loss with fresh mini-batches:

$$L(\theta) = \mathbb{E}_{P \sim \mathcal{P}; (\mathbf{H}; y_{N+1}) \sim P} [\ell_P(\text{read}_y(\text{TF}(\theta); \mathbf{H}); y_{N+1})]; \quad (47)$$

where the loss function $\ell_P : \mathbb{R}^2 \rightarrow \mathbb{R}$ may depend on the training data distribution P in general; we use the square loss when P is regression data, and the logistic loss when P is classification data. We use the Adam optimizer with a fixed learning rate 10^{-4} , which we find works well for all our experiments. Throughout all our experiments except for the sparse linear regression experiment in Figure 3a, we train the model for 300K steps, where each step consists of a (fresh) minibatch with batch size 64 in the base mode, and K minibatches each with batch size 64 in the mixture mode.

For the sparse linear regression experiment, we find that minimizing the training objective (47) alone was not enough, e.g. for the learned transformer to achieve better loss than the least squares algorithm (which achieves much higher test loss than the Lasso; cf. Figure 3a). To help optimization, we augment (47) with another loss that encourages the second-to-last hidden states to recover the true (sparse) coefficient \mathbf{w}^* :

$$L_{\text{t-w}}(\theta) = \frac{1}{N_0} \sum_{j=1}^{N_0} \mathbb{E}_{P \sim \mathcal{P}; (\mathbf{H}; y_{N+1}) \sim P} \left\| \text{TF}^{(1:L-1)}(\theta)_{j:(D-d+1):D} - \mathbf{w}^* \right\|_2^2; \quad (48)$$

Specifically, the above loss encourages the first $N_0 \ll N$ tokens within the second-to-last layer to be close to \mathbf{w}^* . We choose $N_0 = 5$ (recall that the total number of tokens is $N = 10$ and sequence length is $N + 1 = 11$ for this experiment). We minimize the loss $L(\theta) + \lambda L_{\text{t-w}}(\theta)$ with $\lambda = 0.1$ for 2M steps for this task.

Evaluation All evaluations are done on the trained transformer with 6400 test instances. We use the square loss for regression tasks, and the classification error ($1 - \text{accuracy}$) between the true label $y_{N+1} \in \{0, 1\}$ and the predicted label $\hat{y}_{N+1} \in \{0, 1\}$. We report the means in all experiments, as well as their standard deviations (using one-std error bars) in Figure 2a, 2b, 5a, 5b. In Figure 2c, 3b, 3c, 5c, all standard deviations are sufficiently small (not significantly exceeding the width of the markers), thus we did not show error bars in those plots.

Baseline algorithms We implement various baseline machine learning algorithms to compare with the learned transformers. A superset of the algorithms is shown in Figure 3a:

- **Least squares, Logistic regression:** Standard algorithms for linear regression and linear classification, respectively. Note that least squares is also a valid algorithm for classification.
- **Averaging:** The simple algorithm which computes the linear predictor $\hat{\mathbf{w}} = \frac{1}{N} \sum_{i=1}^N y_i \mathbf{x}_i$ and predicts $\hat{y}_{N+1} = \text{sign}(\hat{\mathbf{w}}; \mathbf{x}_{N+1})$;
- **3-NN:** 3-Nearest Neighbors.
- **Ridge:** Standard ridge regression as in (ICRidge). We specifically consider two λ 's (denoted as $\lambda_{1,2}$): $\lambda_1, \lambda_2 = (0.005; 0.125)$. These are the Bayes-optimal regularization strengths for the noise levels $(\sigma_1, \sigma_2) = (0.1; 0.5)$ respectively under the noisy linear model (cf. Corollary 6), using the formula $\lambda^* = d^2 \sigma^2 / N$, with $(d; N) = (20; 40)$.
- **Lasso:** Standard Lasso as in (ICLasso) with $\lambda = (1; 0.1; 0.01; 0.001)g$.

In Figure 2c, the `ridge_analytical` curve plots the expected risk of ridge regression under the noisy linear model over 20 geometrically spaced values of λ 's in between (λ_1, λ_2) , using analytical formulae (with Monte Carlo simulations). The `Bayes_err_{1,2}` indicate the expected risks of λ_1 on task 1 (with noise σ_1) and λ_2 on task 2 (with noise σ_2), respectively.

M.2 Decoder-based architecture

ICL capabilities have also been demonstrated in the literature for decoder-based architectures [31, 2, 47]. There, the transformer can do in-context predictions at every token \mathbf{x}_i using past tokens $\{(\mathbf{x}_j; \mathbf{y}_j)_{j \leq i-1}\}$ as training examples. Here we show that such architectures is also able to perform in-context algorithm selection *at every token*; For results for this architecture on “base” ICL tasks (such as those considered in Figure 3a), we refer the readers to Garg et al. [31].

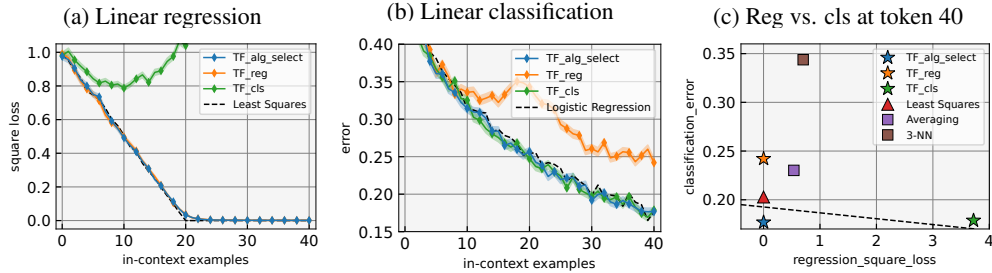


Figure 5: In-context algorithm selection abilities of transformers between linear regression and linear classification. (a,b) On these two tasks, a **single transformer** `TF_alg_select` **simultaneously approaches the performance of the strongest baseline algorithm** Least Squares on linear regression and Logistic Regression on linear classification. (c) At token 40 (using example $f0:::39g$ for training), `TF_alg_select` matches the performance of the best baseline algorithm for both tasks. (a,b,c) Note that transformers pretrained on a single task (`TF_reg`, `TF_cls`) perform near-optimally on their pretraining task but suboptimally on the other task.

Setup Our setup is the same as the two “mixture” modes (linear model + linear classification model, and noisy linear models with two different noise levels) as in Section 6, except that the architecture is GPT-2 following Garg et al. [31], and the input format is changed to (11) (so that the input sequence has $2N + 1$ tokens) without positional encodings. For every $i \geq [N + 1]$, we extract the prediction \hat{y}_i using a linear read-out function applied on output token $2i - 1$, and the (learnable) linear read-out function is the same across all tokens, similar as in Appendix M.1. The rest of the setup (optimization, training, and evaluation) is the same as in Section 6 & M.1. Note that we also train on the objective (47) for all tokens averaged, instead of for the last test token as in Section 6.

Result Figure 2 shows the results for noisy linear models with two different noise levels, and Figure 5 shows the results for linear model + linear classification model. We observe that at every token, In both cases, `TF_alg_select` nearly matches the strongest baseline for both tasks simultaneously, whereas transformers trained on a single task perform suboptimally on the other task. Further, this phenomenon consistently shows up at every token. For example, in Figure 2a & 2b, `TF_alg_select` matches ridge regression with the optimal λ on all tokens $i \geq f1:::Ng$ ($N = 40$). In Figure 5a & 5b, `TF_alg_select` matches least squares on the regression task and logistic regression on the classification task on all tokens $i \geq [N]$. This demonstrates the in-context algorithm selection capabilities of standard decoder-based transformer architectures.

M.3 Computational resource

All our experiments are performed on 8 Nvidia Tesla A100 GPUs (40GB memory). The total GPU time is approximately 5 days (on 8 GPUs), with the largest individual training run taking about a single day on a single GPU.