

## Project Report: Transformers As Statisticians

Ben Agro

## 1 Introduction

Part of the success of large language models (implemented with transformers [1]) can be attributed to their apparent ability to learn new tasks from in-context examples, as opposed to traditional models, which require explicit training for a specific task. Towards a formal understanding of ICL, in this report, we distill two main ideas from the recent paper “Transformers as Statisticians: Provable In-Context Learning with In-Context Algorithm Selection” (TaS) [2]. Section 2 provides preliminaries on transformers and in-context learning. **Main idea 1:** Section 3 discusses how transformers can perform *in-context gradient descent*, which is the underlying mechanism by which transformers can implement the many in-context algorithms, including ridge regression, generalized linear models, and lasso. **Main idea 2:** Section 4 describes one mechanism by which transformers can select the appropriate in-context learning algorithm based on its context. TaS owes its namesake to this mechanism because it is similar to how a statistician can look at data and choose the appropriate model.

## 2 Transformers and In-Context Learning

### Transformers:<sup>1</sup>

Denote the transformer input sequence as a matrix  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_n] \in \mathbf{R}^{D \times N}$ , where each  $\mathbf{h}_i$  is a *token* in the sequence. We will focus on “encoder-mode” transformers, which have multiple layers, and each consisting of a *self-attention* (SA) layer and a residual *multi-layer perceptron* (MLP) layer. Differing from the popular transformer architecture, TaS replaces the softmax used in SA with a (normalized) ReLU activation function  $\sigma(t) = \text{ReLU}(t) = \max(0, t)$ ;

$$\tilde{\mathbf{h}}_i = [\text{SA}_{\theta}(\mathbf{H})]_i = \mathbf{h}_i + \sum_{m=1}^M \frac{1}{N} \sum_{j=1}^N \sigma(\langle \mathbf{Q}_m \mathbf{h}_i, \mathbf{K}_m \mathbf{h}_j \rangle) \mathbf{V}_m \mathbf{h}_j, \quad (2.1)$$

where  $M$  is the number of attention heads, and they do not include layer normalization.<sup>2</sup> The residual MLP layer takes the form  $\text{MLP}_{\mathbf{W}_1, \mathbf{W}_2}(\mathbf{H}) = \mathbf{H} + \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{H})$ , where  $\mathbf{W}_1, \mathbf{W}_2$  are learned weight matrices.

**In-Context Learning:** An ICL instance consists of tokens of the form

$$\mathbf{H} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_N & \mathbf{x}_{N+1} \\ y_1 & y_2 & \dots & y_N & 0 \\ \mathbf{p}_1 & \mathbf{p}_2 & \dots & \mathbf{p}_N & \mathbf{p}_{N+1} \end{bmatrix} \in \mathbf{R}^{D \times (N+1)}, \quad \mathbf{p}_i = \begin{bmatrix} \mathbf{0}_{D-(d+3)} \\ 1 \\ \mathbb{I}\{i < N+1\} \end{bmatrix}, \quad (2.2)$$

where  $\mathbf{x}_i \in \mathbf{R}^d \stackrel{\text{iid}}{\sim} \mathbf{P}_{\mathbf{x}}$  are the input features and  $y_i \in \mathbf{R} \sim \mathbf{P}_{y|\mathbf{x}_i}$  are the labels. Here  $t_i = \mathbb{I}\{i < N+1\}$  indicates the “training”<sup>3</sup> examples; the test example is at token index  $N+1$ . Each ICL instance may have a different data distribution, and the goal is to learn/construct a fixed transformer that can reliably predict the hidden test label  $y_{N+1}$  on large set of data distributions. Note that we include the  $\mathbf{0}_{D-(d+3)}$  in the input token ( $D$  to be specified) to serve as a *placeholder* for the transformer to save intermediate results. Moreover, we can use the 1 to implement constant offsets.

<sup>1</sup>We assume the reader is already familiar with transformers [1].

<sup>2</sup>These changes allow for convenient constructions of transformers in the proofs of TaS, but they are a minor detail as prior work has shown an equivalence between normalized ReLU and softmax activations [3, 4].

<sup>3</sup>In ICL, the model weights are not updated with this data; the prediction is made in the inference pass of the transformer.

### 3 In-Context Learning via Gradient Descent

TaS proves that transformers can implement a broad class of standard machine learning algorithms in context, including least-squares, ridge regression, generalized linear models, and lasso. The linchpin of these results is the construction of a transformer that can perform *in-context gradient descent* (GD).

**Desiderata:** Let  $\hat{L}_N(\mathbf{w}) := \frac{1}{N} \sum_{j=1}^N \ell(\mathbf{w}^T \mathbf{x}_j, y_j)$  be the empirical risk to minimize via GD:

$$\mathbf{w}_{GD}^{t+1} := \mathbf{w}_{GD}^t - \eta \nabla \hat{L}_N(\mathbf{w}_{GD}^t), \quad (3.1)$$

for some loss function  $\ell$ . We want to construct a transformer with  $L + 1$  layers that can approximately implement GD, meaning that  $\mathbf{h}_i^l, \forall l \in \{1, \dots, L\}$  will contain weights  $\hat{\mathbf{w}}^l$  that approximate the GD weights  $\mathbf{w}_{GD}^l$  on the training examples, and the final tokens will contain the predicted labels  $\hat{y}_i = \langle \hat{\mathbf{w}}^L, \mathbf{x}_i \rangle$ . Further, the  $\|\hat{\mathbf{w}}^l - \mathbf{w}_{GD}^l\|$  increases linearly in  $t$ , and so does the test error  $|\hat{y}_{N+1} - y_{N+1}|$  in  $L$ .

**Transformer Construction:** The idea is that the first  $L$  layers will each approximate one step of GD, and that the  $(L + 1)$ -st layer will output the predicted label. Let us start by constructing a single GD step with an attention layer. We assume that the previous layers have an approximation of the weights  $\hat{\mathbf{w}}^t$  (or zero for the initial tokens). We need to assume that  $\partial_s \ell(s, t)$  is *approximable by sum of  $M$  ReLUs*, which means that there are some  $c_m, a_m, b_m, d_m, \forall m \in \{1, \dots, M\}$  such that

$$\sup_{(s,t)} \left| \sum_{m=1}^M c_m \sigma(a_m s + b_m t + d_m) - \partial_s \ell(s, t) \right| \leq \epsilon, \quad (3.2)$$

which is true for a broad class of functions [2]. Note that  $\nabla \hat{L}_N(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N \partial_s \ell(\mathbf{w}^T \mathbf{x}_j, y_j)^T \mathbf{x}_j$ , so that the gradient of the empirical risk and gradient update step based on the training data can be written as:

$$\frac{1}{N} \sum_{m=1}^M \sum_{j=1}^N c_m \sigma(a_m \langle \mathbf{w}, \mathbf{x}_j \rangle + b_m y_j + d_m) \mathbf{x}_j t_j = \nabla \hat{L}_N(\mathbf{w}) + O(\epsilon) \quad (3.3)$$

$$\hat{\mathbf{w}}^l - \frac{1}{N} \sum_{m=1}^M \sum_{j=1}^N \eta c_m \sigma(a_m \langle \hat{\mathbf{w}}^l, \mathbf{x}_j \rangle + b_m y_j + d_m) \mathbf{x}_j t_j = \hat{\mathbf{w}}^l - \eta \nabla \hat{L}_N(\hat{\mathbf{w}}^l) + O(\epsilon). \quad (3.4)$$

We notice the similarity between the LHS of the above expression and the self-attention operation eq. (2.1). Indeed, simple algebra (Proposition C.2 in TaS) shows that for input tokens of the form in eq. (2.2) with the  $\hat{\mathbf{w}}^l$  stored as an intermediate output in the placeholder slot, one can choose  $\mathbf{Q}_m, \mathbf{K}_m, \mathbf{V}_m$  to obtain an SA update that equals the LHS of eq. (3.4). A natural question is whether chaining these approximate GD updates together results in a good approximation of the final weights. One can show (Lemma 14 in TaS) that if we additionally assume  $\ell$  is convex in the first argument, then the weight-space error is linear in  $l$ :  $\|\hat{\mathbf{w}}^l - \mathbf{w}_{GD}^l\| = O(l\epsilon)$ . Thus, we can stack  $L$  of these attention layers (setting the residual MLPs to output identity) to build a transformer that can approximate  $L$  steps of GD with weight-space error in  $O(L\epsilon)$ . Additionally, using the fact that  $\sigma(t) - \sigma(-t) = t$  for  $t > 0$ , it is easy to construct a 2-head attention layer that implements  $\langle \hat{\mathbf{w}}^L, \mathbf{x}_i \rangle$  (each head implements one of the  $\sigma$ 's), and with the linear error bound on  $\hat{\mathbf{w}}^L$ , the prediction error is also linear in  $L$ .

### 4 In-Context Algorithm Selection

Just as a statistician can look at data and choose the appropriate model, TaS shows that transformers can implement *in-context algorithm selection*. We will investigate one of the two proposed mechanisms proposed in for in-context algorithm selection; *post-ICL validation*

**Desiderata:** Assume we have a validation loss  $\hat{L}_{val}(f) = \frac{1}{N_{val}} \sum_{i=1}^{|\mathcal{D}_{val}|} \ell(f(\mathbf{x}_i), y_i)$ ,  $\ell$  is approximable by a sum of ReLUs. We split the in-context examples into a training set  $\mathcal{D}_{train}$  (labelled  $t_i = 1$ ) and a validation set  $\mathcal{D}_{val}$  (labelled  $t_i = -1$ ). The goal is to build a transformer that (1) *Inference*: implements  $K$  in-context algorithms on  $\mathcal{D}_{train}$ ,  $\{f_1, \dots, f_K\}$  and outputs their predictions on all in-context examples, (2) *Evaluation*: approximates the loss  $\tilde{L}_{val}(f_k) \approx \hat{L}_{val}(f_k)$  of each  $f_k$  on  $\mathcal{D}_{val}$ , (3) *Selection*: for each  $x_i$

including the test input ( $t_i = 0$ ), outputs a convex combination  $\hat{f}$  of the predictions  $\{f_1(x_i), \dots, f_K(x_i)\}$  that achieves near the minimum validation loss.

**Transformer construction :** To implement (1), we note that transformers can be *combined in parallel*, meaning that if we have  $K$  transformers (possibly with a different number of layers or heads), each implementing a sequence-to-sequence function  $P_k : \mathbb{R}^{D \times (N+1)} \rightarrow \mathbb{R}^{D \times (N+1)}$ ,

$$\text{TF}_{\theta_k} : \mathbf{H}_k = \begin{bmatrix} \mathbf{H}^{(0)} \\ \mathbf{H}_k^{(1)} \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{H}^{(0)} \\ P_k(\mathbf{H}_k) \end{bmatrix}, \forall k \in \{1, \dots, K\}, \quad (4.1)$$

then, there is a transformer that implements all  $K$  sequence-to-sequence functions in parallel:

$$\text{TF}_{\theta} : \begin{bmatrix} \mathbf{H}^{(0)} \\ \mathbf{H}_1^{(1)} \\ \vdots \\ \mathbf{H}_K^{(1)} \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{H}^{(0)} \\ P_1(\mathbf{H}_1) \\ \vdots \\ P_K(\mathbf{H}_K) \end{bmatrix} \quad (4.2)$$

In section 3, we discussed that transformers can implement many in-context algorithms, so we can stack  $K$  in parallel to implement a single transformer that outputs  $[\ast, f_1(\mathbf{x}_i), \dots, f_K(\mathbf{x}_i), \mathbf{0}_{K+1}]$  in the placeholder slots of  $\mathbf{h}_i$ , where we use  $\ast$  to hide the inputs and labels.

To implement (2), we construct a single attention layer that ingests  $[\ast, f_1(\mathbf{x}_i), \dots, f_K(\mathbf{x}_i), \mathbf{0}_{K+1}]$  and outputs  $[\ast, f_1(\mathbf{x}_i), \dots, f_K(\mathbf{x}_i), \tilde{L}_{val}(f_1), \dots, \tilde{L}_{val}(f_K), 0]$  in the placeholder slots of all  $\mathbf{h}_i$ . Mirroring our discussion of gradient descent in section 3, because  $\ell$  is approximable by a sum of ReLUs, we can construct a single attention layer that implements  $\tilde{L}_{val} \approx \hat{L}_{val}$  using the validation indicator  $t_i = -1$ .

Finally, to implement (3), we construct a three-layer *selection transformer* mapping

$$[\ast, f_1(\mathbf{x}_i), \dots, f_K(\mathbf{x}_i), \tilde{L}_{val}(f_1), \dots, \tilde{L}_{val}(f_K), 0] \rightarrow [\ast, f_1(\mathbf{x}_i), \dots, f_K(\mathbf{x}_i), \tilde{L}_{val}(f_1), \dots, \tilde{L}_{val}(f_K), \hat{f}(\mathbf{x}_i)], \quad (4.3)$$

where  $\hat{f}$  is a convex combination of the  $f_k$  that are within  $\gamma$  of the minimum validation loss:

$$\hat{f} = \sum_{k=1}^K \lambda_k f_k, \text{ s.t } \lambda_1, \dots, \lambda_K \geq 0, \sum_{k=1}^K \lambda_k = 1, \text{ and } \lambda_k > 0 \text{ only if } \tilde{L}_{val}(f_k) \leq \min_{k^* \in [K]} \tilde{L}_{val}(f_{k^*}) + \gamma. \quad (4.4)$$

Define  $c_k = \sum_{p \neq k} \sigma(\tilde{L}_{val}(f_k) - \tilde{L}_{val}(f_p))$ , and notice that if  $c_k \leq \gamma$  then  $\tilde{L}_{val}(f_k) \leq \min_{k^* \in [K]} \tilde{L}_{val}(f_{k^*}) + \gamma$ . We construct the first layer of the selection transformer with a residual MLP (the attention layer implements identity) that maps from the LHS of eq. (4.3) to  $[f_1(\mathbf{x}_1), \dots, f_K(\mathbf{x}_i), c_1, \dots, c_K, 0]$  which is straightforward to construct because  $c_k$  is a sum of ReLUs.

Next, let  $u_k = \sigma(1 - \gamma^{-1}c_k)$ , which is in  $[0, 1]$  and  $u_k > 0$  iff  $c_k \leq \gamma$ . We construct the second layer with an MLP mapping from the above intermediate tokens to  $[f_1(\mathbf{x}_1), \dots, f_K(\mathbf{x}_i), u_1, \dots, u_K, 0]$

Finally, selecting  $\lambda_1 = 1 - \sigma(1 - u_1)$ ,  $\lambda_k = \sigma(1 - u_1 - \dots - u_{k-1}) - \sigma(1 - u_1 - \dots - u_k)$ ,  $\forall k \geq 2$  satisfies the conditions of eq. (4.4). We can re-arrange to get the following expression for  $\hat{f}$  (letting  $f_0 = f_{K+1} = 0$ ):

$$\hat{f} = \sum_{k=1}^{K+1} \sigma(1 - u_1 - \dots - u_{k-1})(f_k(\mathbf{x}_i) - f_{k-1}(\mathbf{x}_i)) = \sum_{k=1}^{K+1} \sum_{w \in \{0,1\}} \sigma(1 - u_1 - \dots - u_{k-1})(-2w + 1)f_{k-w}(\mathbf{x}_i). \quad (4.5)$$

Notice that this has a very similar form to the SA operation in eq. (2.1), so we can implement this with a single attention layer with  $2(K+1)$  heads. Thus, the transformer will predict  $\hat{f}(\mathbf{x}_{N+1})$  on the test example.

Stacking the (1) *inference*, (2) *evaluation*, and (3) *selection* layers gives us a transformer that is a “statistician”, as it can dynamically select the appropriate in-context algorithm based on the context.

## 5 Conclusion

One could imagine how these simple mechanisms of in-context gradient descent and algorithm selection, operating in the high-dimensional representation space of a large transformer, could result in the observed and impressive emergent in-context learning behaviors of LLMs. However, while the TaS’s empirical experiments show that a transformer *can* do these in-context tasks, they do not investigate if it is by the proposed mechanisms, leaving room for future work to understand if large language models are implementing this algorithm selection mechanism in practice.

## References

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [2] Y. Bai, F. Chen, H. Wang, C. Xiong, and S. Mei, “Transformers as statisticians: Provable in-context learning with in-context algorithm selection,” *Advances in neural information processing systems*, vol. 36, 2024.
- [3] M. Wortsman, J. Lee, J. Gilmer, and S. Kornblith, “Replacing softmax with relu in vision transformers,” *arXiv preprint arXiv:2309.08586*, 2023.
- [4] 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.