
Attention with Markov: A Curious Case of Single-layer Transformers

Ashok Vardhan Makkuva^{*1} Marco Bondaschi^{*1} Adway Girish¹ Alliot Nagle² Martin Jaggi¹ Hyeji Kim^{†2}
Michael Gastpar^{†1}

Abstract

In recent years, attention-based transformers have achieved tremendous success across a variety of disciplines including natural languages. To better understand the sequential modeling capabilities of transformers, there is a growing interest in using Markov input processes to study them. While previous research has shown that transformers with two or more layers develop an induction head mechanism to estimate the bigram conditional distribution, we find a surprising empirical phenomenon that single-layer transformers can get stuck at local minima, corresponding to unigrams. To explain this, we introduce a new framework for a principled theoretical and empirical analysis of transformers via Markov chains. Leveraging our framework, we theoretically characterize the loss landscape of single-layer transformers and show the existence of global minima (bigram) and bad local minima (unigram) contingent upon the specific data characteristics and the transformer architecture. Further, we precisely characterize the regimes under which these local optima occur. Backed by experiments, we demonstrate that our theoretical findings are in congruence with the empirical results. Finally, we outline several open problems in this arena. Code is available at <https://anonymous.4open.science/r/Attention-with-Markov-A617/>.

1. Introduction

Attention-based transformers have been at the forefront of recent breakthroughs in a variety of disciplines, including natural language processing (Vaswani et al., 2017; Radford

^{*}Equal contribution ¹School of Computer and Communication Sciences, EPFL, Lausanne, Switzerland ²Department of Electrical and Computer Engineering, UT Austin, Austin, TX, USA. Correspondence to: Ashok Vardhan Makkuva <ashok.makkuva@epfl.ch>.

Work presented at MI workshop at ICML 2024, Vienna, Austria. PMLR 235, 2024. Copyright 2024 by the author(s).

and Narasimhan, 2018; Devlin et al., 2018). One of the key workhorses behind this success is the attention mechanism, which allows transformers to capture complex causal structures in the data, thus rendering them with impressive sequential modeling capabilities.

Given their success, there is tremendous interest and active research in understanding the sequential modeling abilities of transformers. Notably, there is a growing focus on using Markov input processes to study transformers, especially their in-context learning capabilities (Nichani et al., 2024; Edelman et al., 2024; Bietti et al., 2023). These studies reveal an interesting insight that, when the input is a first-order Markov chain, transformers with two or more layers develop an induction head mechanism to estimate the in-context bigram conditional distribution. However, we observe a surprising empirical phenomenon for single-layer transformers, contrary to these findings: contingent on the Markov state switching probabilities, the transformer can get stuck at local minima, corresponding to the unigram rather than the bigram (Fig. 1). This observation, in view of our current limited understanding of transformer models, thus prompts a fundamental question: *Can we systematically characterize the learning capabilities of single-layer transformers with Markovian inputs?*

To address this, in this paper we introduce a new framework for a principled theoretical and empirical analysis of transformers via Markov chains. Leveraging our framework, we characterize the loss landscape of single-layer transformers and prove the existence of bad local minima and global minima corresponding to the unigram and bigram, respectively. Further, we show that the presence of these local optima is contingent upon the Markov state switching probabilities and the weight-tying of the transformer, and precisely characterize the regimes under which this occurs. Together, our analysis reveals an interesting interplay between the data-distributional properties, the transformer architecture, and the final model performance for single-layer transformers with Markov chains.

In summary, we make the following contributions:

- We provide a novel framework for a precise theoretical and empirical study of transformers via Markov chains

(Sec. 3).

- We characterize the loss landscape of single-layer transformers with first-order Markov chains, highlighting the effect of the data distribution and the model architecture (Sec. 4).
- We show that the Markov switching probabilities and weight-tying play a crucial role in the presence of local optima on loss surface and precisely characterize the said conditions (Thms. 2 and 3).

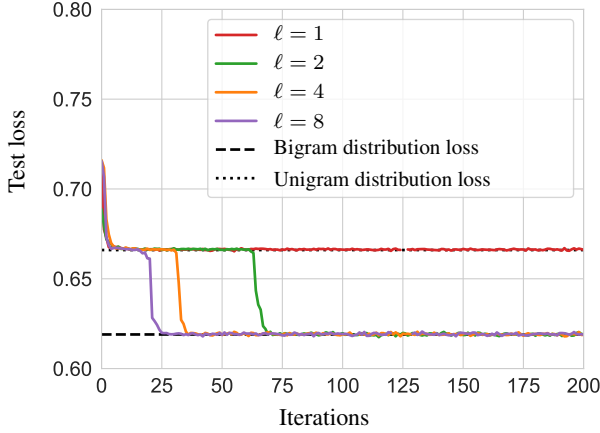


Figure 1: Single-layer transformers ($\ell = 1$) get stuck at local minima, corresponding to the unigram, when the input is a first-order Markov chain with switching probabilities $p = 0.5$ and $q = 0.8$ (Fig. 2b). However, deeper models escape to global minima corresponding to the bigram.

Our main findings and observations are:

- We prove that weight tying can introduce bad local minima for single-layer transformers when Markovian switching is greater than one (Thm. 2). Removing the tying, however, solves the issue (Thm. 3).
- When the Markovian switching is less than one, we empirically observe the model always converges to the global minima irrespective of the weight tying (Fig. 3).
- Interestingly, transformers with depth two and beyond always converge to the global minima irrespective of the weight tying and switching (Fig. 1).

Notation. Scalars are denoted by italic lower case letters like x, y and Euclidean vectors and matrices are denoted by bold ones $\mathbf{x}, \mathbf{y}, \mathbf{M}$, etc. We use $\|\cdot\|$ to denote the ℓ_2 -norm for Euclidean vectors and Frobenius norm for matrices. $[k] \triangleq \{1, \dots, k\}$, and for a sequence $(x_n)_{n \geq 1}$, define $x_k^m \triangleq (x_k, \dots, x_m)$ if $k \geq 1$ and (x_1, \dots, x_m) otherwise. For $z \in \mathbb{R}$, the sigmoid $\sigma(z) \triangleq 1/(1 + e^{-z})$ and

$\text{ReLU}(z) \triangleq \max(0, z)$. For events A and B , $\mathbb{P}(A)$ denotes the probability of A whereas $\mathbb{P}(A | B)$ the conditional probability. Let (x, y) be a pair of discrete random variables on $[k] \times [k]$ with the probability mass function (pmf) of x being $\mathbf{p}_x = (p_1, \dots, p_k) \in [0, 1]^k$. Then its Shannon entropy is defined as $H(x) = H(\mathbf{p}_x) \triangleq -\sum_{i \in [k]} p_i \log p_i$, and the conditional entropy $H(y|x) \triangleq H(x, y) - H(x)$. The entropy rate of a stochastic process $(x_n)_{n \geq 1}$ is defined as $\lim_{n \rightarrow \infty} H(x_1^n)/n$. Finally, for $p \in (0, 1)$, the binary entropy function $h(\cdot)$ is defined as $h(p) \triangleq H(p, 1-p) = -p \log p - (1-p) \log(1-p)$.

2. Background

We describe the transformer architecture and the Markovian input process.

2.1. Transformers

For the ease of exposition, we detail a single-layer transformer with a single-head attention, ReLU non-linearity, and input vocabulary \mathcal{X} of size 2, i.e. $\mathcal{X} = \{0, 1\}$. We omit the layer norm since its influence is marginal in the settings we consider (Sec. 4). Let $\{x_n\}_{n=1}^N \in \{0, 1\}^N$ be an input sequence of length N . Then for each $n \in [N]$, the transformer operations are mathematically given by (Fig. 2a):

$$\mathbf{x}_n = x_n \mathbf{e}_1 + (1 - x_n) \mathbf{e}_0 + \tilde{\mathbf{p}}_n \in \mathbb{R}^d, \quad (\text{Embedding})$$

$$\mathbf{y}_n = \mathbf{x}_n + \mathbf{W}_O \sum_{i \in [n]} \text{att}_{n,i} \cdot \mathbf{W}_V \mathbf{x}_i \in \mathbb{R}^d, \quad (\text{Attention})$$

$$\mathbf{z}_n = \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) \in \mathbb{R}^d, \quad (\text{FF})$$

$$\text{logit}_n = \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}, \quad (\text{Linear})$$

$$f_{\bar{\theta}}(x_1^n) \triangleq \mathbb{P}_{\bar{\theta}}(x_{n+1} = 1 | x_1^n) = \sigma(\text{logit}_n) \in (0, 1). \quad (\text{Prediction})$$

Here $\bar{\theta} \triangleq (\mathbf{e}_1, \mathbf{e}_0, \{\tilde{\mathbf{p}}_n\}, \dots, b, \mathbf{a})$ denotes the full list of the transformer parameters. d is the embedding dimension, \mathbf{e}_1 and \mathbf{e}_0 in \mathbb{R}^d are the token-embeddings corresponding to $x_n = 1$ and $x_n = 0$ respectively, and $\tilde{\mathbf{p}}_n$ is the (trainable) positional encoding. We have matrices $\mathbf{W}_O \in \mathbb{R}^{d \times m}$ and $\mathbf{W}_V \in \mathbb{R}^{m \times d}$, and the attention weights $\text{att}_{n,i} \in (0, 1)$ are computed using the query and key matrices (§ A). $\mathbf{W}_2 \in \mathbb{R}^{d \times r}$ and $\mathbf{W}_1 \in \mathbb{R}^{r \times d}$ are the weight matrices in the FF layer, whereas $\mathbf{a} \in \mathbb{R}^d$ and $b \in \mathbb{R}$ are the weight and bias parameters for the linear layers. For a multi-layer transformer, we apply the successive attention and feed-forward layers multiple times before the final linear layer. Finally, we compute the probability for the symbol 1 using the logits: $f_{\bar{\theta}}(x_1^n) \triangleq \mathbb{P}_{\bar{\theta}}(x_{n+1} = 1 | x_1^n) = \sigma(\text{logit}_n) \in [0, 1]$. Note that a single symbol probability suffices as the vocabulary is binary.

Loss. The parameters $\bar{\theta}$ are trained using the next-token prediction loss between the predicted probability $f_{\bar{\theta}}(x_1^n)$

and the corresponding ground truth symbol x_{n+1} across all the positions $n \in [N]$:

$$L(\bar{\theta}) \triangleq -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [x_{n+1} \cdot \log f_{\bar{\theta}}(x_1^n) + (1 - x_{n+1}) \cdot \log(1 - f_{\bar{\theta}}(x_1^n))], \quad (1)$$

where the expectation is over the data distribution of the sequence $\{x_n\}_{n=1}^N$. In practice, it is replaced by the empirical averages across the sequences $\{x_n\}_{n=1}^N$ sampled from the corpus, with stochastic optimizers like SGD or Adam (Kingma and Ba, 2015) used to update the model parameters.

2.2. Markov chains

We model the input as a *first-order Markov chain*, i.e. a Markov chain with (order) memory $m = 1$. For these processes, the next state is influenced only by the current state and none of the past:

$$\begin{aligned} P_{ij} &\triangleq \mathbb{P}(x_{n+1} = j \mid x_n = i) \\ &= \mathbb{P}(x_{n+1} = j \mid x_n = i, x_1^{n-1} = i_1^{n-1}), \end{aligned}$$

for any $i_1, \dots, i_{n-1}, i, j \in \mathcal{X}, n \geq 1$. Here the Markov kernel $\mathbf{P} = (P_{ij})$ governs the transition dynamics of the process: if $\pi^{(n)} \in [0, 1]^{|\mathcal{X}|}$ denotes the probability law of x_n at time n , then $\pi^{(n+1)} = \pi^{(n)} \cdot \mathbf{P}$. Of particular interest to us in this paper is the kernel $\mathbf{P}(p, q) \triangleq [1 - p, p; q, 1 - q]$ on the binary state space with the switching probabilities $P_{01} = p$ and $P_{10} = q$, for $p, q \in (0, 1)$. Fig. 2b illustrates the state transition diagram for this kernel. Here we refer to the sum $p + q$ as the *switching factor*. We denote a first-order binary Markov chain $(x_n)_{n \geq 1}$ with the transition kernel $\mathbf{P}(p, q)$ and starting with an initial law $\pi^{(1)}$ as $(x_n)_{n \geq 1} \sim (\pi^{(1)}, \mathbf{P}(p, q))$. When the initial distribution is understood from context, we simply write $(x_{n+1} \mid x_n)_{n \geq 1} \sim \mathbf{P}(p, q)$. For this process, the entropy rate equals $H(x_{n+1} \mid x_n) = \frac{1}{p+q} (q h(p) + p h(q))$, which is independent of n .

Stationary distribution. A *stationary distribution* of a Markov chain is a distribution π on \mathcal{X} that is invariant to the transition dynamics, i.e. if $\pi^{(n)} = \pi$, then we have $\pi^{(n+1)} = \pi \mathbf{P} = \pi$ and consequently, $\pi^{(m)} = \pi$ for all $m \geq n$. Also referred to as the steady-state distribution, its existence and uniqueness can be guaranteed under fairly general conditions (Norris, 1997), and in particular for $\mathbf{P}(p, q)$ when $p, q \neq 0, 1$. For $\mathbf{P}(p, q)$, the stationary distribution is given by $\pi(p, q) \triangleq (\pi_0, \pi_1) = \frac{1}{p+q} (q, p)$. The higher the flipping probability q , the higher the likelihood for the chain to be in the state 0 at the steady state. Similarly for the state 1. We can verify that π indeed satisfies $\pi \mathbf{P} = \pi$. For brevity, we drop the dependence on (p, q) and simply write π and \mathbf{P} when the parameters are clear from context.

3. Framework: Transformers via Markov chains

We present our mathematical framework for a principled analysis of transformers via Markov chains. In this paper we focus on first-order Markovian data and single-layer transformers though our framework readily generalizes to higher orders and deeper architectures (Sec. 4.3).

Data. We assume that the vocabulary $\mathcal{X} = \{0, 1\}$ and the input data $\{x_n\}_{n=1}^N \sim (\pi(p, q), \mathbf{P}(p, q))$, for some fixed sequence length $N \geq 1$ and $(p, q) \in (0, 1)^2$. Recall that $p+q$ is the switching factor. The parameters p and q provide a tractable mechanism to control the input data, which plays a crucial role in transformer learning.

Model. We consider a single-layer transformer with a single-head attention, without layer norm. As the input is binary, the [Embedding](#) layer can be simplified to

$$\mathbf{x}_n = x_n \mathbf{e} + \mathbf{p}_n, \quad (\text{Uni-embedding})$$

where $\mathbf{e} \triangleq \mathbf{e}_1 - \mathbf{e}_0$ is the embedding vector and $\mathbf{p}_n \triangleq \mathbf{e}_0 + \tilde{\mathbf{p}}_n$ is the new positional encoding. Note that $x_n \in \{0, 1\}$ and hence the embedding is either $\mathbf{e} + \mathbf{p}_n$ or just \mathbf{p}_n depending on x_n . The other layers are the same as in Sec. 2.1:

$$\begin{aligned} x_n \in \{0, 1\} &\xrightarrow{\text{Uni-embedding}} \mathbf{x}_n \xrightarrow{\text{Attention}} \mathbf{y}_n, \\ \mathbf{y}_n &\xrightarrow{\text{FF}} \mathbf{z}_n \xrightarrow{\text{Linear}} \text{logit}_n \xrightarrow{\text{Prediction}} f_{\bar{\theta}}(x_1^n). \end{aligned} \quad (2)$$

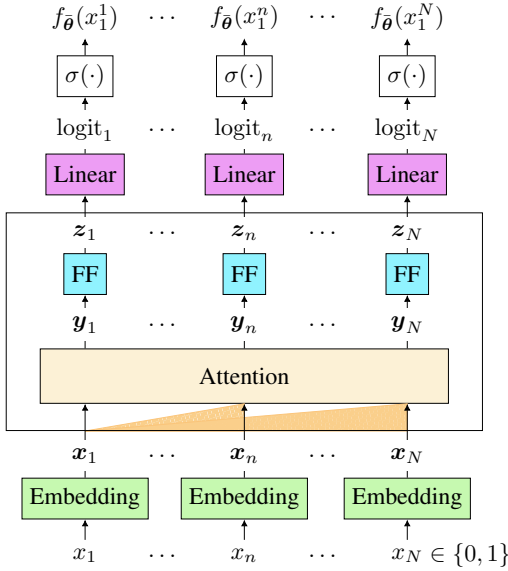
Let $\bar{\theta} \triangleq (\mathbf{e}, \{\mathbf{p}_n\}_{n=1}^N, \dots, b, \mathbf{a}) \in \mathbb{R}^D$ denote the joint list of the parameters from all the layers, with D being the total dimensionality. In training large language models, it is a common practice to tie the embedding and linear layer weights, i.e. $\mathbf{a} = \mathbf{e}$, referred to as *weight tying* (Press and Wolf, 2017). In this case, the list of all parameters, $\theta = (\mathbf{e} = \mathbf{a}, \{\mathbf{p}_n\}_{n=1}^N, \dots, b) \in \mathbb{R}^{D-d}$, since \mathbf{a} is no longer a free parameter. We analyze both weight-tied and general cases.

Loss. We consider the cross-entropy loss L from Eq. (1).

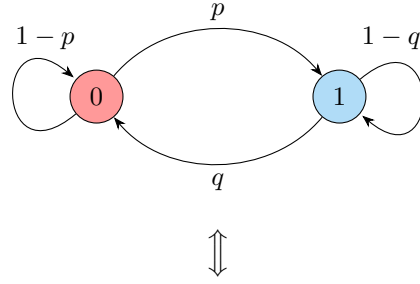
Objective. Towards understanding the phenomenon in Fig. 1, we utilize the aforementioned framework to study single-layer transformers. In particular, our objective is to address the following question:

Can we characterize the loss landscape of single-layer transformers when the input is Markovian?

To build intuition about the loss surface, we first examine its global minima and then provide a detailed characterization of the loss landscape, focusing on local optima, in Sec. 4.



(a) The transformer model with binary input data: for each x_1^n , the next-bit prediction probability is $f_{\theta}(x_1^n) = \mathbb{P}_{\theta}(x_{n+1} = 1 | x_1^n)$.



$$(x_{n+1} | x_n)_{n \geq 1} \sim \mathbf{P}(p, q) \triangleq \begin{bmatrix} 1-p & p \\ q & 1-q \end{bmatrix},$$

$$\mathbf{P}_{ij} = \mathbb{P}(x_{n+1} = j | x_n = i), \quad i, j \in \{0, 1\}.$$

(b) State transition diagram and Markov kernel for a first-order Markov chain $\mathbf{P}(p, q)$ with flipping probabilities $\mathbf{P}_{01} = p$ and $\mathbf{P}_{10} = q$.

Figure 2: Analysis of transformers via Markov chains.

3.1. Single-layer Transformers: Global minima

Since the loss L in Eq. (1) is the cross-entropy loss, it achieves its minimum when the predictive probability matches the Markov kernel (Lemma 1): $f_{\theta}(x_1^n) = \mathbb{P}(x_{n+1} = 1 | x_n)$. In other words, this occurs when the transformer outputs the correct transition probabilities. This raises a natural question: *can a single-layer transformer exactly represent a first-order Markov chain?* Intuitively speaking, this seems plausible since the transformer, even with access to the full past information x_1^n at each $n \in [N]$, can rely solely on the current symbol x_n (Sec. 2). The following result confirms this intuition, showing that such a realization is indeed a *global minimum* for the loss $L(\cdot)$:

Theorem 1 (Global minimum). *Let the input sequence be $\{x_n\}_{n=1}^N \sim (\boldsymbol{\pi}(p, q), \mathbf{P}(p, q))$ for some fixed $(p, q) \in (0, 1)^2$ and $\boldsymbol{\theta} \in \mathbb{R}^{D-d}$ be the transformer parameters for weight-tied case. Then for all (p, q) , there exists a $\boldsymbol{\theta}_{\star} \in \mathbb{R}^{D-d}$ with an explicit construction such that it is a global minimum for the population loss $L(\cdot)$ in Eq. (1) and its prediction matches the Markov kernel, i.e.*

- (i) $L(\boldsymbol{\theta}) \geq L(\boldsymbol{\theta}_{\star})$ for all $\boldsymbol{\theta} \in \mathbb{R}^{D-d}$, and
- (ii) $\mathbb{P}_{\boldsymbol{\theta}_{\star}}(x_{n+1} = 1 | x_1^n) = \mathbb{P}(x_{n+1} = 1 | x_n)$, the Markov kernel or the bigram.

Further, $\boldsymbol{\theta}_{\star}$ satisfies:

- (iii) $L(\boldsymbol{\theta}_{\star}) = H(x_{n+1} | x_n)$, the entropy rate of the Markov chain.
- (iv) $\nabla L(\boldsymbol{\theta}_{\star}) = 0$, i.e. $\boldsymbol{\theta}_{\star}$ is a stationary point.

In addition, the same result holds for the non-weight-tied case when the parameters are in \mathbb{R}^D .

Remark 1. In fact, there exist many such global minima as highlighted in the proof (§ B).

Proof sketch. The key idea here is to show that any $\boldsymbol{\theta}$ satisfying $f_{\theta}(x_1^n) = \mathbb{P}(x_{n+1} = 1 | x_n)$ is a global minimum and is a stationary point with the loss being the entropy rate (Lemmas. 1 and 2). To construct such a $\boldsymbol{\theta}$, we utilize the fact the Markov kernel is only a function of x_n and thus we can ignore the past information in the Attention layer using only the skip. We defer the full proof to § B. \square

Empirical evidence for learning the Markov kernel. As demonstrated in the proof above, a canonical way to realize the Markov kernel by the single-layer transformer is to rely only on the current symbol x_n and ignore the past in the Attention layer. We now empirically confirm this fact. For our experiments, we use the single-layer transformer (Table 1) and report the results averaged across 5 runs and corresponding to the best set of hyper-parameters after a grid search (Table 2). In particular, for $p = 0.2, q = 0.3$, and $d = 4$, we generate sequences $\{x_n\}_{n=1}^N \sim (\boldsymbol{\pi}(p, q), \mathbf{P}(p, q))$ of length $N = 1024$ and train the transformer parameters $\boldsymbol{\theta}$ (weight-tied) to minimize the cross-entropy loss in Eq. (1). At inference, we interpret the attention layer and observe that the relative magnitude of the attention contribution to the final attention output \mathbf{y}_n is negligible, i.e. the ratio $\|\mathbf{W}_O \sum_{i \in [n]} \text{att}_{n,i} \cdot \mathbf{W}_V x_i\| / \|\mathbf{y}_n\| \approx 0.01$. Hence, the attention contribution can be neglected compared to the

skip-connection, i.e. $\mathbf{y}_n \approx \mathbf{x}_n$. Using this approximation, in § D we derive a formula for the final predicted probability $f_{\theta}(x_1^n)$ as it is learnt by the network. This formula reveals interesting insights about the learnt parameters of the transformer:

- **Constant embeddings.** The positional embedding \mathbf{p}_n is constant across n , i.e. it is independent of the sequence position, reflecting the fact that it has learnt to capture the homogeneity of the Markovian chain just from the data.
- **Low-rank weights.** The weight matrices are all approximately rank-one; while it is not fully clear why the training algorithm converges to low rank solutions, they do indeed provide a canonical and simple way to realize the Markov kernel, as illustrated in § D.

Further, we show in § D that plugging in the numerical values obtained from the average of five runs, the probability given by our formula matches the theory, i.e. the model learns to correctly output the Markov kernel probabilities. Indeed, Fig. 3b shows that the test loss of the model converges to the theoretical global minimum (Thm. 1), the entropy rate of the source, corresponding to the bigram, when $p = 0.2$ and $q = 0.3$ ($p + q < 1$). We likewise observe a similar phenomenon without weight tying. For the prediction probability, we focus on the zero positions $n = n_k$ such that $x_{n_k} = 0$. Fig. 3d shows that irrespective of the index k and the past $x_1^{n_k-1}$, if the current bit x_{n_k} is 0, the model correctly predicts the probability for the next bit x_{n_k+1} being 1, which equals p theoretically (Fig. 2b). More precisely, $f_{\theta}(x_1^{n_k-1}, x_{n_k} = 0) = p$ for all $x_1^{n_k-1}$ and k , in line with property (ii) of Thm. 1. A similar conclusion holds with $x_{n_k} = 1$ and prediction probability q . This indicates that the model has learned to recognize the data as first-order Markovian, relying solely on x_n to predict x_{n+1} .

While Thm. 1 and above empirical results highlight the presence of global minima on the loss surface, they does not address local optima, as empirically shown in Fig. 1. We precisely address this in the next section and analyze the loss landscape in terms of the local optima.

4. Single-layer Transformers: Local Optima

In this section we present our main results about the loss landscape of single-layer transformers in terms of local optima. In particular, we prove the existence of bad local minima and saddle points on the loss surface (Thms. 2 and 3), in addition to the global minima discussed above (Thm. 1). Interestingly, the presence of these local optima is influenced by two key factors: *switching factor* of the Markov chain and the *weight tying* of the transformer, highlighting the intricate interplay between the input data and

the model architecture. First, we present the results for the weight tying scenario.

4.1. Weight tying: bad local minima

When the embedding and linear layers are tied, i.e. $\mathbf{e} = \mathbf{a}$, our analysis reveals the following surprising fact: if the switching factor $p + q$ is greater than one, there exist *bad local minima* $\theta_{\pi} \in \mathbb{R}^{D-d}$, where the prediction probability $f_{\theta_{\pi}}(\cdot)$ is the marginal stationary distribution π (unigram), disregarding the past and the present information (Thm. 2 and Fig. 3c). Now we state the result formally. Let $L_{\star} \triangleq L(\theta_{\star})$ denote the global minimal loss from Thm. 1.

Theorem 2 (Bad local minimum). *Let the input sequence be $\{x_n\}_{n=1}^N \sim (\pi(p, q), \mathbf{P}(p, q))$ for a fixed $(p, q) \in (0, 1)^2$ and the transformer parameters be weight-tied. If $p + q > 1$, there exists an explicit $\theta_{\pi} \in \mathbb{R}^{D-d}$ such that it is a bad local minimum for the loss $L(\cdot)$, i.e.*

- there exists a neighborhood $\mathcal{B}(\theta_{\pi}, r)$ with $r > 0$ such that $L(\theta) \geq L(\theta_{\pi})$ for all $\theta \in \mathcal{B}(\theta_{\pi}, r)$, with $L(\theta_{\pi}) > L_{\star}$.*

Further, θ_{π} satisfies:

- $\mathbb{P}_{\theta_{\pi}}(x_{n+1} = 1 | x_1^n) = \mathbb{P}(x_{n+1} = 1) = \pi_1$, the marginal distribution or the unigram.*
- $L(\theta_{\pi}) = H(x_{n+1}) = H(\pi)$, the entropy of the marginal.*
- $\nabla L(\theta_{\pi}) = 0$, i.e. θ_{π} is a stationary point.*

Remark 2. Since $L(\theta_{\pi}) = H(x_{n+1})$ and $L_{\star} = H(x_{n+1}|x_n)$, the optimality gap $L(\theta_{\pi}) - L_{\star} = H(x_{n+1}) - H(x_{n+1}|x_n) = I(x_n; x_{n+1}) \geq 0$, where $I(x_n; x_{n+1})$ is the mutual information between x_n and x_{n+1} (Cover and Thomas, 2006). It equals zero if and only if x_n and x_{n+1} are independent, which happens for $p + q = 1$ (since $\mathbb{P}(x_{n+1} = 1 | x_n) = x_n(1 - p - q) + p$).

Proof sketch. The main idea behind constructing θ_{π} is that if we set $\mathbf{e} = \mathbf{a} = 0$ in the Linear layer, the model ignores the inputs all together and outputs a constant probability, and in particular π_1 by choosing the bias b appropriately, i.e. $f_{\theta_{\pi}}(x_1^n) = \pi_1$ for all x_1^n, n . For this θ_{π} it's easy to show that $L(\theta_{\pi}) = H(\pi)$ and that it's a stationary point. Further we show that the Hessian at θ_{π} follows the structure $\begin{bmatrix} \mathbf{H}_{\alpha} & 0 \\ 0 & 0 \end{bmatrix}$ where $\mathbf{H}_{\alpha} \succ 0$ when $p + q > 1$, and that it implies the local minimality of θ_{π} . We defer the full proof to § B.4. \square

Empirical evidence for bad local minima. The proof of Thm. 2 above highlights that when the linear weight \mathbf{a} is zero in θ , it serves as a bad local minimum. While this might seem as a theoretical anomaly, we empirically confirm that it is not. Specifically, we use the same weight-tied setting

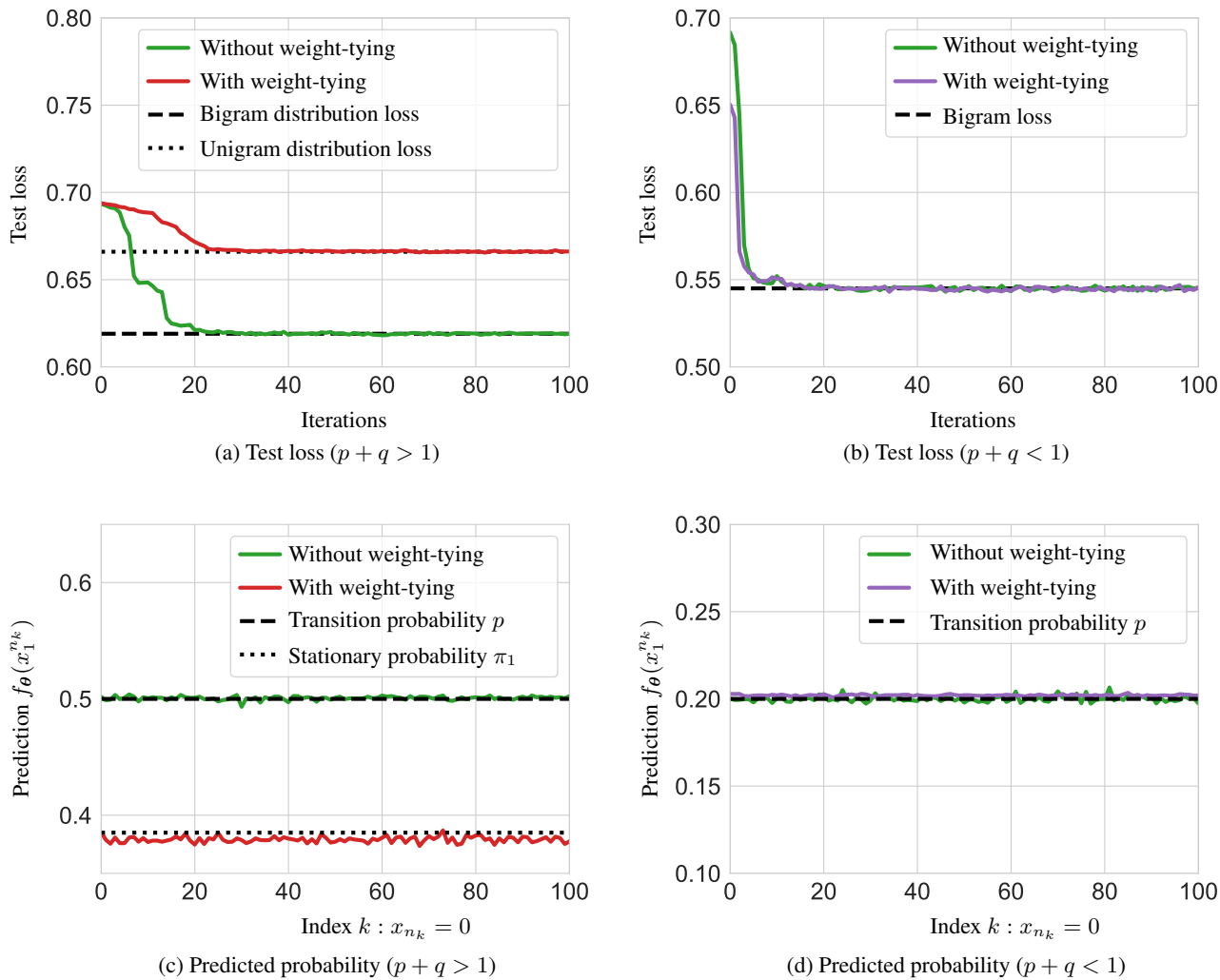


Figure 3: Effect of weight tying on test loss and predicted probabilities $f_{\theta}(x_1^{n_k})$ for zero indices $\{n_k\}_{k=1}^{100}$ such that $x_{n_k} = 0$. For (a),(c): $p = 0.5, q = 0.8$. With weight tying, the loss converges to a local minimum, and the predicted probability is $\pi_1 = p/(p + q)$. Without weight tying, we predict the correct probability p and converge to a global minimum. For (b),(d): $p = 0.2, q = 0.3$. The test loss always converges to a global minimum, and the predicted probability is p .

as before but with the flipping probabilities $p = 0.5$ and $q = 0.8$ instead, we observe that the magnitude of the vector \mathbf{a} is approximately 0.01 whereas that of \mathbf{z}_n is 4.0 in the Linear layer. Thus, $\langle \mathbf{a}, \mathbf{z}_n \rangle \approx 0.04$, while the bias term $b \approx -0.4$ implying $\sigma(\langle \mathbf{a}, \mathbf{z}_n \rangle + b) \approx \sigma(b)$. Hence the final prediction returned by the network only depends on the bias of the linear layer, totally independent of the input sequence x_1^N . Fig. 3c illustrates this fact and shows that the model always predicts the stationary probability π_1 at all positions in the sequence, independent of the input. Here we plot for the indices k such that $x_{n_k} = 0$ but we observe the same phenomenon for $x_{n_k} = 1$, i.e. $f_\theta(x_1^n) = \pi_1$ for all x_1^n, n , thus verifying property (ii) of Thm. 2. Similarly, Fig. 3a demonstrates that the model test loss converges to the entropy of the stationary distribution $H(\pi)$, the unigram loss, instead of the global minimum L_\star given by the entropy rate of the source (Thm. 2 - (iii)). Fig. 4 illustrates for a wider range of $(p, q) \in (0, 1)^2$.

Interpreting global and local minima. Thm. 2 should be interpreted in the light that it guarantees the existence of bad local minima only for $p + q > 1$ (in sync with the experiments, Fig. 4). While there could exist such minima even when $p + q < 1$, we empirically observe that the model always converges to the global minimum in this scenario (Fig. 3b). Likewise, while Thm. 1 guarantees the existence of global minima for all (p, q) and in particular for $p + q > 1$, empirically the model often converges to bad local minima as highlighted above (Fig. 3a).

4.2. Without weight tying: saddle points

Now we let the token-embedding $e \in \mathbb{R}^d$ and the linear weight $\mathbf{a} \in \mathbb{R}^d$ be independent parameters. Interestingly, here, the earlier local minimum θ_π becomes a saddle point.

Theorem 3 (Saddle point). *Consider the same setting as in Thm. 2 and for $p + q > 1$, let $\theta_\pi = (e_\pi = \mathbf{a}_\pi, \dots, b_\pi) \in \mathbb{R}^{D-d}$ be the corresponding bad local minimum for the loss $L(\cdot)$ in the weight-tied scenario. Then its extension $\hat{\theta}_\pi \triangleq (\theta_\pi, \mathbf{a}_\pi) \in \mathbb{R}^D$ is a saddle point for $L(\cdot)$ in \mathbb{R}^D in the non-weight-tied case. Further, $\hat{\theta}_\pi$ satisfies the same properties (ii)–(iv) as in Thm. 2.*

Empirical evidence and interpretation. In view of the theoretical results above, removing weight tying is possibly beneficial: the bad local minimum in the weight-tied case for $p + q > 1$ suddenly becomes a saddle point when the weight tying is removed. We observe a similar phenomenon empirically (Fig. 4): as shown in Fig. 3a, when not weight-tied, the model’s test loss converges to the entropy rate of the source when $p + q > 1$, in contrast to the weight-tied case, possibly escaping the saddle point (Thm. 3). The fact that the model eventually learns the correct Markovian kernel is further demonstrated by the red curve in Fig. 3c. Figs. 3b and 3a together highlight that the model always (empirically) con-

verges to the global minimum in the non-weight-tied case irrespective of the switching factor $p + q$.

Key insights. Together, our theoretical and empirical results highlight that when the switching $p + q > 1$, the weight-tied model can get stuck at bad local minima corresponding to the unigram. In contrast, the non-weight-tied model can potentially escape saddle points to converge to the global minima, corresponding to the bigram (Markov kernel). This explains the phenomenon in Fig. 1 for the single-layer transformer, where $p = 0.5$ and $q = 0.8$. When $p + q < 1$, we empirically observe that the model always converges to a global minimum irrespective of weight-tying.

4.3. Does depth help escape local minima?

For single-layer transformers, the aforementioned results highlight the significance of the switching factor and the weight tying on the loss landscape. In stark contrast, we empirically observe that for transformers of depth 2 and beyond, the loss curve eventually reaches the global minimum regardless of these factors, as highlighted in Fig. 1. Interestingly, we observe during training that it first reaches a plateau at a loss value corresponding to $H(\pi)$ and after a few additional iterations, it further drops down until it reaches the global minimum (the entropy rate). This suggests that, while there could still be local minima, increasing the number of layers positively affects the loss curvature at the bad local minima in a manner making it easier to escape and reach the global minimum. In the context of feed-forward neural networks, depth of the architecture has been shown to play a major role in terms of the representation power and learning capabilities (Telgarsky, 2016). Given our empirical observations, a similar analysis for transformers that demonstrates the benefits of depth is an intriguing direction for future research.

5. Related work

There is tremendous interest and active research in understanding transformer models from various perspectives ((Giannou et al., 2023; Oymak et al., 2023; Li et al., 2023a; Fu et al., 2023; Noci et al., 2023; Tian et al., 2023) and references therein). Yun et al. (2020); Pérez et al. (2021); Wei et al. (2022); Malach (2023); Jiang and Li (2023) demonstrate the representation capabilities of transformers and show properties such as universal approximation and Turing-completeness. Another line of inquiry (Elhage et al., 2021; Snell et al., 2021; Wang et al., 2023) is mechanistic interpretability, i.e. reverse-engineering transformer operations on specific synthetic tasks (e.g., matrix inversion and eigenvalue decomposition in Charton (2022), modular addition in Nanda et al. (2023)) to understand the transformer components but they usually lack theoretical guarantees, as opposed to ours. Li et al. (2023c) studies how trans-

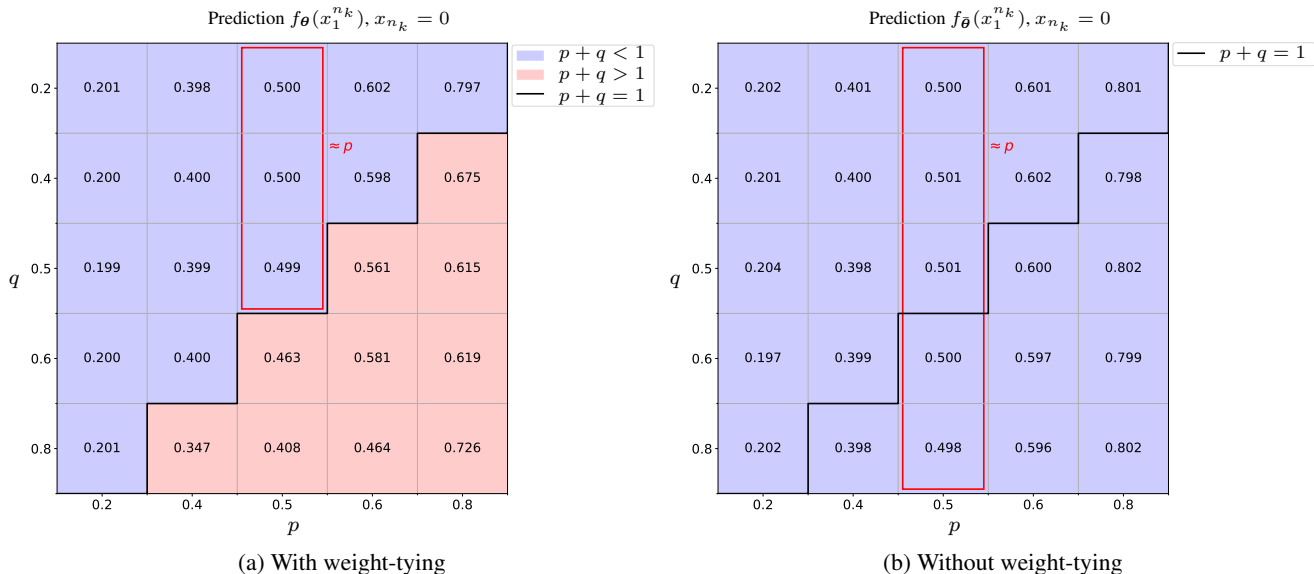


Figure 4: Average of predicted probabilities across 5 runs for different values of p and q , with and without weight-tying. In the former case, there is a clear demarcation between the cases where $p + q < 1$ and those where $p + q > 1$. For $p + q < 1$, all runs accurately predict the correct conditional probability. For $p + q > 1$, some of the runs predict the stationary probability instead, causing the average to diverge from the correct p . In the latter case, the model always predicts the correct probability for all p and q .

formers learn semantic structures across words while we are interested in how they learn sequentiality in input data. (Tarzanagh et al., 2023a;b) take an optimization-theoretic perspective to study training dynamics and characterize implicit biases in transformer models trained with gradient descent. In contrast, we characterize the local and global minima of the loss landscape of these models under sequential input data.

(Chen et al., 2024; Dong et al., 2023; Akyürek et al., 2023; Von Oswald et al., 2023; Xie et al., 2021; Bai et al., 2023; Li et al., 2023b; Garg et al., 2022) study in-context learning, i.e. the ability of the transformer to extract the desired task from just a few representative examples. While we consider the transformer architecture in fully generality including ReLU nonlinearity, (Bietti et al., 2023) assumes frozen position encodings, query matrices, and linear activations, whereas (Edelman et al., 2024) assumes a minimal model for two-layer transformer with only trainable attention component and linear activation to analyze a single gradient-descent step. (Grau-Moya et al., 2024) use data generated from Markov chains, among other data sources, to study if meta-learning can approximate Solomonoff Induction. (Chung et al., 2021) provide empirical evidence to suggest that weight tying has drawbacks in encoder-only models, which is in line with our observations that removing weight tying is beneficial in decoder-only models with Markovian input data. More recently, (Rajaraman et al., 2024) study the effect of tokenization on learning Markov

chains. (Ildiz et al., 2024) show an equivalence between the attention mechanism and Markov models, whereas we characterize the loss landscape of attention-based transformers when the input is Markovian.

6. Conclusion and Open questions

In this work, we provide a novel framework for a systematic theoretical and empirical study of the sequential modeling capabilities of transformers through Markov chains. Leveraging this framework, we theoretically characterize the loss landscape of single-layer transformers and show the existence of global minima and bad local minima contingent upon the specific data characteristics and the transformer architecture, and independently verify them by experiments. We further reveal interesting insights for deeper architectures. We believe our framework provides a new avenue for a principled study of transformers. In particular, some interesting open questions in this direction include:

- Characterizing the learning dynamics of gradient-based algorithms in our setup.
- Understanding the interplay between the depth of the transformer and the order of the Markov process.

References

- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? Investigations with linear models. In *The Eleventh International Conference on Learning Representations*, 2023.
- Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. In *Workshop on Efficient Systems for Foundation Models @ ICML2023*, 2023.
- Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- François Charton. What is my math transformer doing? – Three results on interpretability and generalization, 2022. URL <https://arxiv.org/abs/2211.00170>.
- Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Training dynamics of multi-head softmax attention for in-context learning: Emergence, convergence, and optimality, 2024.
- Hyung Won Chung, Thibault Fevry, Henry Tsai, Melvin Johnson, and Sebastian Ruder. Rethinking embedding coupling in pre-trained language models. In *International Conference on Learning Representations*, 2021.
- Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons, 2nd edition, 2006.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding, 2018. URL <https://arxiv.org/abs/1810.04805>.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. A survey on in-context learning, 2023. URL <https://arxiv.org/abs/2301.00234>.
- Benjamin L. Edelman, Ezra Edelman, Surbhi Goel, Eran Malach, and Nikolaos Tsilivis. The evolution of statistical induction heads: In-context learning markov chains, 2024.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova Das-Sarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. URL <https://transformer-circuits.pub/2021/framework/index.html>.
- Hengyu Fu, Tianyu Guo, Yu Bai, and Song Mei. What can a single attention layer learn? A study through the random features lens. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? A case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- Angeliki Giannou, Shashank Rajput, Jy-Yong Sohn, Kangwook Lee, Jason D. Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers. In *Proceedings of the 40th International Conference on Machine Learning*, pages 11398–11442, 23–29 Jul 2023.
- Jordi Grau-Moya, Tim Genewein, Marcus Hutter, Laurent Orseau, Grégoire Delétang, Elliot Catt, Anian Ruoss, Li Kevin Wenliang, Christopher Mattern, Matthew Aitchison, and Joel Veness. Learning universal predictors, 2024. URL <https://arxiv.org/abs/2401.14953>.
- Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
- M. Emrullah Ildiz, Yixiao Huang, Yingcong Li, Ankit Singh Rawat, and Samet Oymak. From self-attention to markov models: Unveiling the dynamics of generative transformers, 2024.
- Haotian Jiang and Qianxiao Li. Approximation theory of transformer networks for sequence modeling, 2023. URL <https://arxiv.org/abs/2305.18475>.
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity. In *The Eleventh International Conference on Learning Representations*, 2023a.
- Yingcong Li, M. Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: generalization and stability in in-context learning. In *Proceedings of the 40th International Conference on Machine Learning*, 2023b.
- Yuchen Li, Yuanzhi Li, and Andrej Risteski. How do transformers learn topic structure: towards a mechanistic understanding. In *Proceedings of the 40th International Conference on Machine Learning*, 2023c.

- Eran Malach. Auto-regressive next-token predictors are universal learners, 2023.
- Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. In *The Eleventh International Conference on Learning Representations*, 2023.
- Eshaan Nichani, Alex Damian, and Jason D Lee. How transformers learn causal structure with gradient descent. *arXiv preprint arXiv:2402.14735*, 2024.
- Lorenzo Noci, Chuning Li, Mufan Bill Li, Bobby He, Thomas Hofmann, Chris J. Maddison, and Daniel M. Roy. The shaped transformer: Attention models in the infinite depth-and-width limit. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- J. R. Norris. *Markov Chains*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1997.
- Samet Oymak, Ankit Singh Rawat, Mahdi Soltanolkotabi, and Christos Thrampoulidis. On the role of attention in prompt-tuning. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.
- Matteo Pagliardini. GPT-2 modular codebase implementation. <https://github.com/epfnl/llm-baselines>, 2023.
- Ofir Press and Lior Wolf. Using the output embedding to improve language models. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 157–163, April 2017. URL <https://aclanthology.org/E17-2025>.
- Jorge Pérez, Pablo Barceló, and Javier Marinkovic. Attention is Turing-complete. *Journal of Machine Learning Research*, 22(75):1–35, 2021.
- Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-training. 2018. URL <https://api.semanticscholar.org/CorpusID:49313245>.
- Nived Rajaraman, Jiantao Jiao, and Kannan Ramchandran. Toward a theory of tokenization in LLMs, 2024.
- Charlie Snell, Ruiqi Zhong, Dan Klein, and Jacob Steinhardt. Approximating how single head attention learns, 2021. URL <https://arxiv.org/abs/2103.07601>.
- Davoud Ataee Tarzanagh, Yingcong Li, Christos Thrampoulidis, and Samet Oymak. Transformers as support vector machines. In *NeurIPS 2023 Workshop on Mathematics of Modern Machine Learning*, 2023a.
- Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Max-margin token selection in attention mechanism. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023b.
- Matus Telgarsky. Benefits of depth in neural networks. In *Conference on learning theory*, pages 1517–1539. PMLR, 2016.
- Yuandong Tian, Yiping Wang, Beidi Chen, and Simon Shaolei Du. Scan and snap: Understanding training dynamics and token composition in 1-layer transformer. In *Conference on Parsimony and Learning (Recent Spotlight Track)*, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.
- Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pages 35151–35174, 2023.
- Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The Eleventh International Conference on Learning Representations*, 2023.
- Colin Wei, Yining Chen, and Tengyu Ma. Statistically meaningful approximation: a case study on approximating Turing machines with transformers. In *Advances in Neural Information Processing Systems*, volume 35, pages 12071–12083, 2022.
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit Bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021. URL <https://arxiv.org/abs/2111.02080>.
- Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. Are transformers universal approximators of sequence-to-sequence functions? In *International Conference on Learning Representations*, 2020.

Organization. The appendix is organized as follows:

- App. A details the transformer architecture, especially that of the attention mechanism.
- App. B provides the proofs for our theoretical results on first-order Markov chains.
- App. C contains additional experimental details and results for first-order Markov chains.

A. The Transformer architecture

We describe the Transformer architecture from Sec. 2.1 in detail, using the embedding layer simplification from Sec. 3:

$$\begin{aligned}
 \mathbf{x}_n &= x_n \mathbf{e} + \mathbf{p}_n \in \mathbb{R}^d, && \text{(Uni-embedding)} \\
 \mathbf{y}_n &= \mathbf{x}_n + \mathbf{W}_O \sum_{i \in [n]} \text{att}_{n,i} \cdot \mathbf{W}_V \mathbf{x}_i \in \mathbb{R}^d, && \text{(Attention)} \\
 \mathbf{z}_n &= \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) \in \mathbb{R}^d, && \text{(FF)} \\
 \text{logit}_n &= \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}, && \text{(Linear)} \\
 f_{\bar{\theta}}(x_1^n) &\triangleq \mathbb{P}_{\bar{\theta}}(x_{n+1} = 1 \mid x_1^n) = \underbrace{\sigma(\text{logit}_n)}_{\in [0,1]}. && \text{(Prediction)}
 \end{aligned}$$

(i) Embedding: The discrete tokens $x_n = 1$ and $x_n = 0$ are mapped to the token-embeddings \mathbf{e} and 0 in \mathbb{R}^d respectively, where d is the embedding dimension. The positional embedding $\mathbf{p}_n \in \mathbb{R}^d$ encodes the positional information (varies with n). The sum of these two embeddings constitutes the final input embedding $\mathbf{x}_n \in \mathbb{R}^d$.

(ii) Attention: The attention layer can be viewed as mapping a query and a set of key-value pairs to an output, which are all vectors (Vaswani et al., 2017). That is, on top of the skip-connection \mathbf{x}_n , the output $\mathbf{y}_n \in \mathbb{R}^d$ is computed as a weighted sum of the values $\mathbf{v}_i \triangleq \mathbf{W}_V \mathbf{x}_i$. The weight assigned to each value, $\text{att}_{n,i}$, is computed by a compatibility function of the query vector $\mathbf{q}_n \triangleq \mathbf{W}_Q \mathbf{x}_n$ and the corresponding key vectors $\mathbf{k}_i \triangleq \mathbf{W}_K \mathbf{x}_i$ for all $i \in [n]$. More precisely, $\text{att}_{n,i} \triangleq \text{softmax}(\langle \mathbf{q}_n, \mathbf{k}_1 \rangle, \dots, \langle \mathbf{q}_n, \mathbf{k}_n \rangle) / \sqrt{d}$. $\mathbf{W}_{K,Q,V} \in \mathbb{R}^{m \times d}$ are the respective key, query, and value matrices, and $\mathbf{W}_O \in \mathbb{R}^{d \times m}$ is the projection matrix. For multi-headed attention, the same operation is performed on multiple parallel heads, whose outputs are additively combined.

(iii) Feed-forward (FF): The FF transformation consists of a skip-connection and a single-hidden layer with ReLU activation and weight matrices $\mathbf{W}_2 \in \mathbb{R}^{d \times r}$, and $\mathbf{W}_1 \in \mathbb{R}^{r \times d}$. The FF layer is applied token-wise on each $\mathbf{y}_n \in \mathbb{R}^d$ to output $\mathbf{z}_n \in \mathbb{R}^d$ with the same dimensionality.

(iv) Linear: The linear layer transforms the final output embedding \mathbf{z}_n to a scalar $\text{logit}_n \in \mathbb{R}$, with the weight parameter $\mathbf{a} \in \mathbb{R}^d$ and the bias $b \in \mathbb{R}$.

(v) Prediction: The sigmoid activation finally converts the scalar logits to probabilities for the next-token prediction. Since the vocabulary has only two symbols, it suffices to compute the probability for the symbol 1: $f_{\bar{\theta}}(x_1^n) \triangleq \mathbb{P}_{\bar{\theta}}(x_{n+1} = 1 \mid x_1^n) = \sigma(\text{logit}_n) \in [0, 1]$. More generally, the logits are of the same dimensionality as the vocabulary and are converted to the prediction probabilities using a softmax layer, which simplifies to the sigmoid for the binary case. Likewise, there are as many token-embeddings as the words in the vocabulary and several layers of multi-headed attention and FF operations are applied alternatively on the input embeddings to compute the final logits.

Finally, The Transformer parameters $\bar{\theta} \triangleq (\mathbf{e}, \{\mathbf{p}_n\}_{n=1}^N, \dots, b, \mathbf{a}) \in \mathbb{R}^D$ are trained via the cross-entropy loss on the next-token prediction:

$$L(\bar{\theta}) \triangleq -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [x_{n+1} \cdot \log f_{\bar{\theta}}(x_1^n) + (1 - x_{n+1}) \cdot \log(1 - f_{\bar{\theta}}(x_1^n))], \quad (3)$$

B. Proofs of Sec. 4

We now present our proofs for the technical results in Sec. 4. Towards this, we first establish two useful lemmas on the loss function $L(\cdot)$ and the corresponding gradient computation. Let $\bar{\theta} = (e, \{\mathbf{p}_n\}_{n=1}^N, \dots, b, \mathbf{a}) \in \mathbb{R}^D$ be the list of parameters in the non-weight-tied case and $\theta = (e = \mathbf{a}, \{\mathbf{p}_n\}_{n=1}^N, \dots, b) \in \mathbb{R}^{D-d}$ in the weight-tied case. With a slight abuse of notation, by $w \in \bar{\theta}$ we mean a specific parameter w among the set $\{e, \mathbf{p}_1, \dots, \mathbf{p}_N, \dots, b, \mathbf{a}\}$. Since the weight-tied scenario is a special case of the non-weight-tied one with $e = \mathbf{a}$, we directly present the results for the general non-weight-tied case, but both lemmas hold for $\theta \in \mathbb{R}^{D-d}$ as well. First we start with the loss function.

Lemma 1 (Loss as KL divergence). *Let the input sequence be $\{x_n\}_{n=1}^N \sim (\pi(p, q), \mathbf{P}(p, q))$ for some fixed $(p, q) \in (0, 1)^2$, $\bar{\theta} = (e, \{\mathbf{p}_n\}_{n=1}^N, \dots, b, \mathbf{a}) \in \mathbb{R}^D$ be the full list of the transformer parameters, and $L(\bar{\theta})$ be the corresponding cross-entropy loss in Eq. (1). Then the loss function $L(\cdot)$ is equivalent to the KL divergence between the Markov kernel and the predicted distribution, i.e.*

$$L(\bar{\theta}) = \frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} [D_{\text{KL}}(\mathbb{P}(x_{n+1} = \cdot | x_n) \parallel \mathbb{P}_{\bar{\theta}}(x_{n+1} = \cdot | x_1^n))] + H(x_{n+1} | x_n), \quad (4)$$

where $D_{\text{KL}}(P \parallel Q)$ is the KL divergence between two distributions P and Q , and $H(x_{n+1} | x_n)$ is the entropy rate of the Markov chain.

Remark 3. Consequently, Eq. (4) highlights that any parameter $\bar{\theta}$ with the predicted probability $f_{\bar{\theta}}(x_1^n) = \mathbb{P}(x_{n+1} = 1 | x_n)$ is a global minimum for the loss L as $D_{\text{KL}}(\cdot \parallel \cdot) \geq 0$ (Cover and Thomas, 2006). We utilize this fact in the proof of Thm. 1 below.

Proof. We defer the proof to § B.6. □

Lemma 2 (Gradient computation). *Consider the same data and parameter setting as in Lemma 1 and $L(\bar{\theta})$ be the cross-entropy loss in Eq. (1). Then for any parameter $w \in \bar{\theta}$,*

$$\begin{aligned} \nabla_w L(\bar{\theta}) &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [(x_{n+1} - f_{\bar{\theta}}(x_1^n)) \cdot \nabla_w (\mathbf{a}^\top \mathbf{z}_n + b)] \\ &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} [(\mathbb{P}(x_{n+1} = 1 | x_n) - f_{\bar{\theta}}(x_1^n)) \cdot \nabla_w (\mathbf{a}^\top \mathbf{z}_n + b)]. \end{aligned} \quad (5)$$

Remark 4. Eq. (5) highlights that any parameter $\bar{\theta}$ with the predicted probability $f_{\bar{\theta}}(x_1^n) = \mathbb{P}(x_{n+1} = 1 | x_n)$ is also a stationary point for the loss L . We utilize this fact in the proof of Thm. 1 below.

Proof. We defer the proof to § B.7. □

We now detail the proofs of theorems in Sec. 4. We prove the global minimum result in Thm. 1 in two parts, separately for the cases when $p + q \leq 1$ and $p + q > 1$. For both these cases, first we consider weight-tying and in App. B.3, the non-tied case.

B.1. Proof of Thm. 1 for $p + q \leq 1$, weight-tying

Proof. We assume that $p + q \leq 1$ and that we use weight tying, i.e. the list of parameters $\theta = (e = \mathbf{a}, \{\mathbf{p}_n\}_{n=1}^N, \dots, b) \in \mathbb{R}^{D-d}$. Thus in view of Lemma 1 and Lemma 2, it follows that any θ satisfying $f_{\theta}(x_1^n) = \mathbb{P}(x_{n+1} = 1 | x_n)$ is a global minimum with loss equalling the entropy rate, and is a stationary point. Hence it suffices to construct such a θ .

To build our intuition towards designing θ_* , recall that the Markov kernel $\mathbb{P}(x_{n+1} = 1 | x_n)$ can be succinctly written as $\mathbb{P}(x_{n+1} = 1 | x_n) = x_n(1 - p - q) + p$. To ensure that $f_{\theta}(x_1^n) = x_n(1 - p - q) + p$, it suffices for the transformer to utilize only the information from the current symbol x_n and ignore the past x_1^{n-1} . In view of the transformer architecture in § A, a natural way to realize this is to let $\mathbf{W}_O = 0$ and $\mathbf{W}_2 = 0$ in Attention and FF respectively. This implies that $\mathbf{z}_n = \mathbf{y}_n = \mathbf{x}_n = x_n \mathbf{e} + \mathbf{p}_n$. Hence the logits are given by $\text{logit}_n = \langle \mathbf{e}, \mathbf{z}_n \rangle + b = x_n \|\mathbf{e}\|^2 + \langle \mathbf{e}, \mathbf{p}_n \rangle + b$. Since $f_{\theta}(x_1^n) = \sigma(\text{logit}_n)$ and it equals the Markov kernel, we have that

$$\sigma(\text{logit}_n) = \sigma(x_n \|\mathbf{e}\|^2 + \langle \mathbf{e}, \mathbf{p}_n \rangle + b) = x_n(1 - p - q) + p.$$

Rewriting,

$$x_n \|e\|^2 + \langle e, \mathbf{p}_n \rangle + b = \log \left(\frac{x_n(1-p-q) + p}{(1-p)(1-x_n) + qx_n} \right), \quad x_n \in \{0, 1\}.$$

Substituting $x_n = 0$ and $x_n = 1$, we further simplify to

$$\begin{aligned} \langle e, \mathbf{p}_n \rangle + b &= \log \left(\frac{p}{1-p} \right), \\ \|e\|^2 + \langle e, \mathbf{p}_n \rangle + b &= \log \left(\frac{1-q}{q} \right). \end{aligned}$$

Subtracting both the equations we obtain that a global minimum θ should satisfy

$$\begin{aligned} \|e\|^2 &= \log \left(\frac{(1-p)(1-q)}{pq} \right), \\ \langle e, \mathbf{p}_n \rangle + b &= \log \left(\frac{p}{1-p} \right). \end{aligned} \tag{6}$$

Note the above choice of e is well-defined since $\frac{(1-p)(1-q)}{pq} > 1$ when $p+q < 1$ and hence $\log \frac{(1-p)(1-q)}{pq} > 0$. While there exist infinitely many solutions for (e, \mathbf{p}_n, b) satisfying Eq. (6), a canonical such solution for the global minimum $\theta = \theta_*$ is

$$\theta_* = \left(e = \mathbf{a} = \mathbf{1} \sqrt{\frac{1}{d} \log \frac{(1-p)(1-q)}{pq}}, \{\mathbf{p}_n = 0\}_{n=1}^N, \mathbf{W}_O = 0, \mathbf{W}_{K,Q,V}, \mathbf{W}_2 = 0, \mathbf{W}_1, b = \log \frac{p}{1-p} \right), \tag{7}$$

where $\mathbf{1} \in \mathbb{R}^{D-d}$ denotes the all-one vector, the position embeddings \mathbf{p}_n are set to zero, $\mathbf{W}_{K,Q,V} \in \mathbb{R}^{m \times d}$, and $\mathbf{W}_1 \in \mathbb{R}^{r \times d}$ can be set to any arbitrary value. This concludes the explicit construction of θ_* and the proof. \square

B.2. Proof of Thm. 1 for $p+q > 1$, weight-tying

Proof. We use a similar idea as in the proof for the $p+q \leq 1$ case by constructing a $\theta \in \mathbb{R}^{D-d}$ satisfying $f_\theta(x_1^n) = \mathbb{P}(x_{n+1} = 1 | x_n) = x_n(1-p-q) + p$. However, in this case we need to use the ReLU component of the FF mechanism unlike the earlier case where we set $\mathbf{W}_2 = 0$. Now we start with constructing θ_* .

Let the embedding $e = \mathbf{a} = \mathbf{1}$ and the positional encoding $\mathbf{p}_n = -\frac{1}{2}\mathbf{1}$ for all $n \geq 1$, where $\mathbf{1} \in \mathbb{R}^d$ denotes the all-one vector. Thus $x_n = \alpha_n \mathbf{1}$ with $\alpha_n = +\frac{1}{2}$ when $x_n = 1$ and $\alpha_n = -\frac{1}{2}$ when $x_n = 0$. Now let $\mathbf{W}_O = 0$ in the Attention layer. Hence $\mathbf{y}_n = x_n = \alpha_n \mathbf{1}$. For the FF layer, let \mathbf{W}_1 and \mathbf{W}_2 be such that (to be determined later)

$$\mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) = \beta_n \mathbf{1},$$

and hence

$$\mathbf{z}_n = x_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) = \alpha_n \mathbf{1} + \beta_n \mathbf{1} = (\alpha_n + \beta_n) \mathbf{1}.$$

Thus the logits are given by $\text{logit}_n = \sigma(\langle \mathbf{a}, \mathbf{z}_n \rangle + b) = \sigma(d(\alpha_n + \beta_n) + b)$. Since $f_\theta(x_1^n) = \sigma(\text{logit}_n) = \mathbb{P}(x_{n+1} = 1 | x_n)$, we have that

$$\sigma(\text{logit}_n) = \sigma(d(\alpha_n + \beta_n) + b) = x_n(1-p-q) + p, \quad x_n \in \{0, 1\}.$$

Rewriting,

$$d(\alpha_n + \beta_n) + b = \log \left(\frac{x_n(1-p-q) + p}{(1-p)(1-x_n) + qx_n} \right), \quad x_n \in \{0, 1\}.$$

Substituting $x_n = 0$ and $x_n = 1$, and denoting corresponding β 's by β_1 and β_0 (with a slight abuse of notation), we further simplify to

$$d \left(-\frac{1}{2} + \beta_0 \right) + b = \log \left(\frac{p}{1-p} \right), \tag{8}$$

$$d \left(\frac{1}{2} + \beta_1 \right) + b = \log \left(\frac{1-q}{q} \right).$$

Subtracting both the above equations we obtain

$$d(1 + \beta_1 - \beta_0) = \underbrace{\log \left(\frac{(1-p)(1-q)}{pq} \right)}_{<0 \text{ when } p+q>1}. \quad (9)$$

Now it suffices to find β_1 and β_0 satisfying Eq. (9). Recall that β_n obeys $\mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) = \beta_n \mathbf{1}$. Let $\mathbf{W}_1 = w \mathbf{1} \mathbf{1}^\top$ and $\mathbf{W}_2 = -\mathbf{W}_1^\top$ for some $w \in \mathbb{R}$. Since $\mathbf{y}_n = \alpha_n \mathbf{1}$, we have

$$-w \mathbf{1} \mathbf{1}^\top \text{ReLU}(w \mathbf{1} \mathbf{1}^\top \alpha_n \mathbf{1}) = \beta_n \mathbf{1}.$$

Simplifying,

$$\beta_n = -w^2 d \cdot \mathbf{1}^\top \text{ReLU}(\alpha_n \mathbf{1}), \quad \alpha_n = \pm \frac{1}{2}.$$

Thus $\beta_0 = 0$ (corresponding to $x_n = 0$ and $\alpha_n = -\frac{1}{2}$) and $\beta_1 = -\frac{w^2 d^2}{2}$ (otherwise). Substituting them in Eq. (9), we have

$$1 - \frac{w^2 d^2}{2} = \frac{1}{d} \cdot \log \left(\frac{(1-p)(1-q)}{pq} \right).$$

Let $w = w_*$ be a solution to the above equation, i.e.

$$w_* = \sqrt{\frac{2}{d^2} \left(1 - \frac{1}{d} \cdot \log \left(\frac{(1-p)(1-q)}{pq} \right) \right)}.$$

By substituting $\beta_0 = 0$ in Eq. (8) we obtain the bias $b_* = \log \left(\frac{p}{1-p} \right) + \frac{d}{2}$. Piecing everything together, let

$$\boldsymbol{\theta}_* = \left(\mathbf{e} = \mathbf{a} = \mathbf{1}, \{\mathbf{p}_n = (-1/2) \mathbf{1}\}_{n=1}^N, \mathbf{W}_O = 0, \mathbf{W}_{K,Q,V}, \mathbf{W}_1 = w_* \mathbf{1} \mathbf{1}^\top, \mathbf{W}_2 = -\mathbf{W}_1^\top, b = b_* \right), \quad (10)$$

and we are done. \square

B.3. Proof of Thm. 1 for non-weight-tied

Proof. By extending $\boldsymbol{\theta}_* \in \mathbb{R}^{D-d}$ to $\bar{\boldsymbol{\theta}}_* \triangleq (\boldsymbol{\theta}_*, \mathbf{a}_*) \in \mathbb{R}^D$, it follows from the Transformer architecture in § A that $\mathbb{P}_{\bar{\boldsymbol{\theta}}_*}(x_{n+1} = 1 | x_1^n) = \mathbb{P}_{\boldsymbol{\theta}_*}(x_{n+1} = 1 | x_1^n) = \mathbb{P}(x_{n+1} = 1 | x_n)$, the Markov kernel. As the proof of Thm. 1 in § B.2 and § B.1 establish, prediction probability equalling the kernel is a sufficient condition for global-optimality. Hence $\bar{\boldsymbol{\theta}}_*$ is a global minimum for $L(\cdot)$ in \mathbb{R}^D . \square

B.4. Proof of Thm. 2

Proof. First we construct an explicit $\boldsymbol{\theta}_\pi \in \mathbb{R}^{D-d}$ such that it satisfies properties (ii)–(iv) of Thm. 2 i.e. it is a stationary point with loss value being the entropy of the marginal $H(\pi)$ and that it captures the marginal distribution $\mathbb{P}(x_{n+1} = 1) = \pi_1$. Then we compute its Hessian and show that it is a local minimum for $p + q > 1$ thus proving property (i). On the other hand, the same $\boldsymbol{\theta}_\pi$ could either be a local minimum or saddle point for $p + q < 1$. We start with the construction.

Recall that the full set of the Transformer parameters in the weight-tied case is given by $\boldsymbol{\theta} = (\mathbf{e} = \mathbf{a}, \{\mathbf{p}_n\}_{n=1}^N, \mathbf{W}_O, \mathbf{W}_{K,Q,V}, \mathbf{W}_2, \mathbf{W}_1, b) \in \mathbb{R}^{D-d}$. Define $\boldsymbol{\theta}_\pi \in \mathbb{R}^{D-d}$ to be

$$\boldsymbol{\theta}_\pi = \left(\mathbf{e} = \mathbf{a} = 0, \{\mathbf{p}_n\}_{n=1}^N, \mathbf{W}_O = 0, \mathbf{W}_{K,Q,V}, \mathbf{W}_2 = 0, \mathbf{W}_1, b = \log \left(\frac{p}{q} \right) \right), \quad (11)$$

where $\{\mathbf{p}_n\}_{n=1}^N \subset \mathbb{R}^d$, $\mathbf{W}_{K,Q,V} \in \mathbb{R}^{m \times d}$, and $\mathbf{W}_1 \in \mathbb{R}^{r \times d}$ can be set to any arbitrary value. Now we start with property (ii).

(ii): $f_{\boldsymbol{\theta}_\pi}(x_1^n) = \mathbb{P}_{\boldsymbol{\theta}_\pi}(x_{n+1} = 1 \mid x_1^n) = \mathbb{P}(x_{n+1} = 1) = \pi_1$.

Since $\mathbf{a} = 0$, it follows from (Linear) and (Prediction) layers that $f_{\boldsymbol{\theta}_\pi}(x_1^n) = \sigma(b) = \sigma(\log(p/q)) = \frac{p}{p+q} = \pi_1$. In other words, the model ignores all the inputs and outputs a constant probability π_1 .

(iii): $L(\boldsymbol{\theta}_\pi) = H(x_{n+1}) = H(\boldsymbol{\pi})$.

Since $f_{\boldsymbol{\theta}_\pi}(\cdot) = \pi_1 = \mathbb{E}[x_{n+1}]$, it follows from Eq. (3) that

$$\begin{aligned} L(\boldsymbol{\theta}_\pi) &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [x_{n+1} \cdot \log f_{\boldsymbol{\theta}_\pi}(x_1^n) + (1 - x_{n+1}) \cdot \log(1 - f_{\boldsymbol{\theta}_\pi}(x_1^n))] \\ &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [x_{n+1} \cdot \log \pi_1 + (1 - x_{n+1}) \cdot \log \pi_0] \\ &= \frac{1}{N} \sum_{n \in [N]} [-\pi_1 \log \pi_1 - \pi_0 \log \pi_0] \\ &= H(\boldsymbol{\pi}) = H(x_{n+1}). \end{aligned}$$

(iv): $\nabla L(\boldsymbol{\theta}_\pi) = 0$.

At $\boldsymbol{\theta} = \boldsymbol{\theta}_\pi$, the individual layer outputs of the Transformer (§ A) are given by

$$x_n \in \{0, 1\} \xrightarrow{\text{Uni-embedding}} \mathbf{x}_n = \mathbf{p}_n \xrightarrow{\text{Attention}} \mathbf{y}_n = \mathbf{p}_n \xrightarrow{\text{FF}} \mathbf{z}_n = \mathbf{p}_n \xrightarrow{\text{Linear}} \text{logit}_n = b \xrightarrow{\text{Prediction}} f_{\boldsymbol{\theta}_\pi}(x_1^n) = \pi_1.$$

In other words, none of the layer outputs depend on the input sequence $\{x_n\}_{n=1}^N$. In view of this fact and $\mathbb{E}[x_{n+1}] = \pi_1$, using Lemma 2 the gradient with respect to \mathbf{a} of L at $\boldsymbol{\theta} = \boldsymbol{\theta}_\pi$ is given by

$$\begin{aligned} \nabla_{\mathbf{a}} L &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [(x_{n+1} - f_{\boldsymbol{\theta}_\pi}(x_1^n)) \cdot \nabla_{\mathbf{a}} (\mathbf{a}^\top \mathbf{z}_n + b)] \\ &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [(x_{n+1} - \pi_1) (\mathbf{z}_n + \nabla_{\mathbf{a}} \mathbf{z}_n \cdot \mathbf{a})] \\ &\stackrel{(\mathbf{a}=0)}{=} -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_{n+1}} [(x_{n+1} - \pi_1) \cdot \mathbf{p}_n] \\ &= 0. \end{aligned}$$

Similarly, for b :

$$\nabla_b L = -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [(x_{n+1} - f_{\boldsymbol{\theta}_\pi}(x_1^n)) \cdot \nabla_b (\mathbf{a}^\top \mathbf{z}_n + b)] = -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_{n+1}} [x_{n+1} - \pi_1] = 0.$$

For any other parameter $\mathbf{w} \in \boldsymbol{\theta}$ apart from \mathbf{a}, b , we see from Eq. (5) that the gradient $\nabla_{\mathbf{w}} L$ has the term $\nabla_{\mathbf{w}} (\mathbf{a}^\top \mathbf{z}_n) = (\nabla_{\mathbf{w}} \mathbf{z}_n) \cdot \mathbf{a}$ inside the expectation $\mathbb{E}_{x_1^{n+1}}[\dots]$. Since $\mathbf{a} = 0$, this equals zero and hence $\nabla_{\mathbf{w}} L = 0$.

Together $\nabla L(\boldsymbol{\theta}_\pi) = 0$.

(i): $\boldsymbol{\theta}_\pi$ is a bad local minimum for L when $p + q > 1$.

Towards establishing this, we first let $\boldsymbol{\alpha} = (b, \mathbf{a})$ and $\boldsymbol{\beta} = (\{\mathbf{p}_n\}_{n=1}^N, \mathbf{W}_O, \mathbf{W}_{K,Q,V}, \mathbf{W}_2, \mathbf{W}_1)$ be two different sets of parameters comprising $\boldsymbol{\theta}$, i.e. $\boldsymbol{\theta} = (\boldsymbol{\alpha}, \boldsymbol{\beta})$ and compute the Hessian $\mathbf{H}_\pi \triangleq \nabla^{(2)} L(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_\pi}$ and show that it has the following block-diagonal structure:

$$\mathbf{H}_\pi = \begin{bmatrix} \mathbf{H}_\alpha & 0 \\ 0 & 0 \end{bmatrix},$$

where H_α corresponds to the Hessian with respect to the parameters \mathbf{a} and b in α . Further we show that if $p + q > 1$, $H_\alpha \succ 0$ i.e. it is positive-definite. This helps us in establishing that θ_π a local minimum. Now we start with the Hessian computation.

Hessian computation. We first compute the Hessian with respect to α .

From Lemma 2, we have that second derivative with respect to b at $\theta = \theta_\pi$ is given by

$$\begin{aligned} \nabla_b^{(2)} L &= \nabla_b (\nabla_b L) = \nabla_b \left(-\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [x_{n+1} - f_\theta(x_1^n)] \right) \\ &\stackrel{(a)}{=} \frac{1}{N} \sum_{n \in [N]} \mathbb{E} [f_\theta(x_1^n)(1 - f_\theta(x_1^n)) \cdot \nabla_b (\mathbf{a}^\top \mathbf{z}_n + b)] \\ &\stackrel{(\theta = \theta_\pi)}{=} \frac{1}{N} \sum_{n \in [N]} \mathbb{E} [\pi_1 \pi_0] \\ &= \pi_0 \pi_1 > 0, \end{aligned}$$

where (a) follows from the fact that $\nabla_b f_\theta(x_1^n) = \nabla_b \sigma(\mathbf{a}^\top \mathbf{z}_n + b) = f_\theta(x_1^n)(1 - f_\theta(x_1^n)) \cdot \nabla_b (\mathbf{a}^\top \mathbf{z}_n + b)$. Now we compute the second derivative with respect to \mathbf{a} . From Lemma 2, we obtain

$$\begin{aligned} \nabla_{\mathbf{a}}^{(2)} L &= \nabla_{\mathbf{a}} (\nabla_{\mathbf{a}} L) = \nabla_{\mathbf{a}} \left(-\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [(x_{n+1} - f_\theta(x_1^n)) \cdot \nabla_{\mathbf{a}} (\mathbf{a}^\top \mathbf{z}_n)] \right) \\ &= \nabla_{\mathbf{a}} \left(-\frac{1}{N} \sum_{n \in [N]} \mathbb{E} [(x_{n+1} - f_\theta(x_1^n))(z_n + (\nabla_{\mathbf{a}} \mathbf{z}_n) \cdot \mathbf{a})] \right) \\ &\stackrel{(a)}{=} \frac{1}{N} \sum_{n \in [N]} \mathbb{E} [f_\theta(1 - f_\theta)(z_n + (\nabla_{\mathbf{a}} \mathbf{z}_n) \cdot \mathbf{a})(z_n + (\nabla_{\mathbf{a}} \mathbf{z}_n) \cdot \mathbf{a})^\top] \\ &\quad - \frac{1}{N} \sum_{n \in [N]} \mathbb{E} [(x_{n+1} - f_\theta(x_1^n))(2\nabla_{\mathbf{a}} \mathbf{z}_n)], \end{aligned}$$

where (a) follows from the gradient of the product rule and the fact that $\nabla_{\mathbf{a}} f_\theta(x_1^n) = f_\theta(x_1^n)(1 - f_\theta(x_1^n))(z_n + (\nabla_{\mathbf{a}} \mathbf{z}_n) \cdot \mathbf{a})$. At $\theta = \theta_\pi$, this further simplifies to

$$\begin{aligned} \nabla_{\mathbf{a}}^{(2)} L &\stackrel{(b)}{=} \frac{1}{N} \sum_{n \in [N]} (\mathbb{E} [\pi_1 \pi_0 \cdot \mathbf{p}_n \mathbf{p}_n^\top] - 2\mathbb{E} [(x_{n+1} - \pi_1)x_n \mathbf{I}]) \\ &\stackrel{(c)}{=} \frac{1}{N} \sum_{n \in [N]} ((\pi_0 \pi_1) \cdot \mathbf{p}_n \mathbf{p}_n^\top) - 2\mathbb{E} [(x_n(1 - p - q) + p - \pi_1)x_n \mathbf{I}] \\ &= \frac{1}{N} \sum_{n \in [N]} ((\pi_0 \pi_1) \cdot \mathbf{p}_n \mathbf{p}_n^\top) - 2\mathbb{E} [x_n(\pi_0 - q) \mathbf{I}] \\ &= \frac{1}{N} \sum_{n \in [N]} ((\pi_0 \pi_1) \cdot \mathbf{p}_n \mathbf{p}_n^\top) - 2\pi_1(\pi_0 - q) \mathbf{I} \\ &= \pi_0 \pi_1 \left(\sum_{n \in [N]} \frac{\mathbf{p}_n \mathbf{p}_n^\top}{N} - 2 \left(1 - \frac{q}{\pi_0} \right) \mathbf{I} \right) \\ &\stackrel{(d)}{=} \pi_0 \pi_1 \left(\sum_{n \in [N]} \frac{\mathbf{p}_n \mathbf{p}_n^\top}{N} + 2(p + q - 1) \mathbf{I} \right), \end{aligned}$$

where (b) follows from the fact that $\nabla_{\mathbf{a}} \mathbf{z}_n = x_n \mathbf{I}$ at $\boldsymbol{\theta} = \boldsymbol{\theta}_\pi$ where \mathbf{I} is the identity matrix is $\mathbb{R}^{d \times d}$, (c) from the observation that $\mathbb{E}[x_{n+1}|x_n] = x_n(1-p-1) + p$, and (d) from the fact that $\pi_0 = \frac{q}{p+q}$. Now we compute the cross-derivative of second order $\nabla_{\mathbf{a}, b} L$. Again, invoking Lemma 2,

$$\begin{aligned} \nabla_{\mathbf{a}b} L &= \nabla_{\mathbf{a}}(\nabla_b L) = \nabla_{\mathbf{a}} \left(-\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [x_{n+1} - f_{\boldsymbol{\theta}}(x_1^n)] \right) \\ &= \frac{1}{N} \sum_{n \in [N]} \mathbb{E} [f_{\boldsymbol{\theta}}(1 - f_{\boldsymbol{\theta}})(\mathbf{z}_n + (\nabla_{\mathbf{a}} \mathbf{z}_n) \cdot \mathbf{a})] \\ &\stackrel{(\boldsymbol{\theta} = \boldsymbol{\theta}_\pi)}{=} \frac{1}{N} \sum_{n \in [N]} \mathbb{E} [\pi_1 \pi_0 \cdot \mathbf{p}_n] \\ &= \pi_0 \pi_1 \left(\sum_{n \in [N]} \frac{\mathbf{p}_n}{N} \right). \end{aligned}$$

Piecing all the results together we obtain that for $\boldsymbol{\alpha} = (b, \mathbf{a})$, its corresponding Hessian is given by

$$\mathbf{H}_{\boldsymbol{\alpha}} \triangleq \nabla_{\boldsymbol{\alpha}}^{(2)} L(\boldsymbol{\alpha}_\pi, \boldsymbol{\beta}_\pi) = \pi_0 \pi_1 \begin{bmatrix} 1 & \mathbf{u}^\top \\ \mathbf{u} & \mathbf{V} \end{bmatrix}, \quad \mathbf{u} \triangleq \sum_{n \in [N]} \frac{\mathbf{p}_n}{N}, \mathbf{V} \triangleq \sum_{n \in [N]} \frac{\mathbf{p}_n \mathbf{p}_n^\top}{N} + 2(p+q-1)\mathbf{I}. \quad (12)$$

We now show that the Hessian $\mathbf{H}_{\boldsymbol{\beta}} \triangleq \nabla_{\boldsymbol{\beta}}^{(2)} L(\boldsymbol{\alpha}_\pi, \boldsymbol{\beta}_\pi) = 0$. Recall that $\boldsymbol{\beta} = (\{\mathbf{p}_n\}_{n=1}^N, \mathbf{W}_O, \mathbf{W}_{K,Q,V}, \mathbf{W}_2, \mathbf{W}_1)$. For any $\mathbf{w}_1, \mathbf{w}_2 \in \boldsymbol{\beta}$, Lemma 2 implies that

$$\begin{aligned} \nabla_{\mathbf{w}_1 \mathbf{w}_2} L &= \nabla_{\mathbf{w}_1}(\nabla_{\mathbf{w}_2} L) = \nabla_{\mathbf{w}_1} \left(-\frac{1}{N} \sum_{n \in [N]} \mathbb{E} [(x_{n+1} - f_{\boldsymbol{\theta}}(x_1^n))(\nabla_{\mathbf{w}_2} \mathbf{z}_n \cdot \mathbf{a})] \right) \\ &\stackrel{(a)}{=} \frac{1}{N} \sum_{n \in [N]} \mathbb{E} [f_{\boldsymbol{\theta}}(1 - f_{\boldsymbol{\theta}})(\nabla_{\mathbf{w}_2} \mathbf{z}_n \cdot \mathbf{a})(\nabla_{\mathbf{w}_1} \mathbf{z}_n \cdot \mathbf{a})^\top] \\ &\stackrel{(\boldsymbol{\theta} = \boldsymbol{\theta}_\pi)}{=} 0, \end{aligned}$$

where (a) follows from the fact that $\nabla_{\mathbf{w}_1} f_{\boldsymbol{\theta}}(x_1^n) = \nabla_{\mathbf{w}_1} \sigma(\mathbf{a}^\top \mathbf{z}_n + b) = f_{\boldsymbol{\theta}}(x_1^n)(1 - f_{\boldsymbol{\theta}}(x_1^n))(\nabla_{\mathbf{w}_1} \mathbf{z}_n \cdot \mathbf{a})$. Thus $\mathbf{H}_{\boldsymbol{\beta}} = 0$. Similarly, we can show that $\mathbf{H}_{\boldsymbol{\alpha}\boldsymbol{\beta}} = \nabla_{\boldsymbol{\alpha}\boldsymbol{\beta}} L = 0$ and hence $\mathbf{H}_{\boldsymbol{\beta}\boldsymbol{\alpha}} = \mathbf{H}_{\boldsymbol{\alpha}\boldsymbol{\beta}}^\top = 0$. Thus,

$$\mathbf{H}_{\boldsymbol{\pi}} = \nabla^{(2)} L(\boldsymbol{\theta}_\pi) = \begin{bmatrix} \mathbf{H}_{\boldsymbol{\alpha}} & 0 \\ 0 & 0 \end{bmatrix}$$

Now it remains to show that $\mathbf{H}_{\boldsymbol{\alpha}}$ is positive-definite when $p+q > 1$ and it implies that $\boldsymbol{\theta}_\pi$ is a local minimum.

Positive-definiteness of $\mathbf{H}_{\boldsymbol{\alpha}}$. Recall from Eq. (12) that $\mathbf{H}_{\boldsymbol{\alpha}} = \begin{bmatrix} 1 & \mathbf{u}^\top \\ \mathbf{u} & \mathbf{V} \end{bmatrix}$, where $\mathbf{u} = \sum_{n \in [N]} \mathbf{p}_n / N$, $\mathbf{V} = \sum_{n \in [N]} \mathbf{p}_n \mathbf{p}_n^\top / N + 2(p+q-1)\mathbf{I}$. From the characterization of positive-definiteness by Schur's complement (Horn and Johnson, 2012), we have that $\mathbf{H}_{\boldsymbol{\alpha}} \succ 0 \Leftrightarrow 1 > 0$ and $\mathbf{V} - \mathbf{u}\mathbf{u}^\top \succ 0$. We have that

$$\begin{aligned} \mathbf{V} - \mathbf{u}\mathbf{u}^\top &= 2(p+q-1)\mathbf{I} + \sum_{n \in [N]} \frac{\mathbf{p}_n \mathbf{p}_n^\top}{N} - \left(\sum_{n \in [N]} \frac{\mathbf{p}_n}{N} \right) \left(\sum_{n \in [N]} \frac{\mathbf{p}_n}{N} \right)^\top \\ &= 2(p+q-1)\mathbf{I} + \sum_{n \in [N]} \frac{1}{N} \left(\mathbf{p}_n - \sum_{n \in [N]} \frac{\mathbf{p}_n}{N} \right) \left(\mathbf{p}_n - \sum_{n \in [N]} \frac{\mathbf{p}_n}{N} \right)^\top \\ &= 2(p+q-1)\mathbf{I} + \text{Cov}(\{\mathbf{p}_n\}_{n=1}^N), \end{aligned}$$

where $\text{Cov}(\{\mathbf{p}_n\}_{n=1}^N) = \sum_{n \in [N]} \frac{1}{N} \left(\mathbf{p}_n - \sum_{n \in [N]} \frac{\mathbf{p}_n}{N} \right) \left(\mathbf{p}_n - \sum_{n \in [N]} \frac{\mathbf{p}_n}{N} \right)^\top$ is the covariance matrix of the set $\{\mathbf{p}_n\}_{n=1}^N$ and hence positive semi-definite. Thus if $p + q > 1$, we have that $2(p + q - 1)\mathbf{I} \succ 0$ and together, we obtain that $\mathbf{V} - \mathbf{u}\mathbf{u}^\top \succ 0$. Hence $\mathbf{H}_\alpha \succ 0$. Now it remains to show that $\boldsymbol{\theta}_\pi$ is a local minimum.

\mathbf{H}_α is positive-definite implies $\boldsymbol{\theta}_\pi$ is a local minimum. Since $\mathbf{H}_\alpha \succ 0$, let $\mathbf{H}_\alpha \succcurlyeq \lambda \mathbf{I}$ for some $\lambda > 0$ (in fact $\lambda = 2(p + q - 1)$ works). Since $\boldsymbol{\theta} = (\boldsymbol{\alpha}, \boldsymbol{\beta}) \in \mathbb{R}^{D-d}$, interpret $L(\boldsymbol{\theta}) = L(\boldsymbol{\alpha}, \boldsymbol{\beta})$ as a function of two variables $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ with appropriate dimensions. We know the following facts about $L(\cdot, \cdot)$:

- **Fact 1.** $\boldsymbol{\alpha} \mapsto L(\boldsymbol{\alpha}, \boldsymbol{\beta}_\pi)$ has a local minimum (as a function of one variable) at $\boldsymbol{\alpha} = \boldsymbol{\alpha}_\pi$ (since $\mathbf{H}_\alpha \succ 0$).
- **Fact 2.** $\boldsymbol{\beta} \mapsto L(\boldsymbol{\alpha}_\pi, \boldsymbol{\beta})$ is constant in $\boldsymbol{\beta}$ (since $\mathbf{a}_\pi = 0$, the probability $f_\theta(x_1^n)$ is constant w.r.t. z_n and hence w.r.t. $\boldsymbol{\beta}$ (Linear)).
- **Fact 3.** $\nabla L(\boldsymbol{\alpha}_\pi, \boldsymbol{\beta}_\pi) = 0$ and $\mathbf{H}_\pi = \nabla^{(2)} L(\boldsymbol{\alpha}_\pi, \boldsymbol{\beta}_\pi) = \begin{bmatrix} \mathbf{H}_\alpha & 0 \\ 0 & 0 \end{bmatrix}$ with $\mathbf{H}_\alpha \succcurlyeq \lambda \mathbf{I}$.

Using these facts now we show that $(\boldsymbol{\alpha}_\pi, \boldsymbol{\beta}_\pi) = \boldsymbol{\theta}_\pi$ is also a local minimum in two-variables. We prove this by contradiction. Suppose that $(\boldsymbol{\alpha}_\pi, \boldsymbol{\beta}_\pi)$ is not a local minimum for $L(\cdot, \cdot)$. Without loss of generality, by a shift of coordinates treat $\boldsymbol{\theta}_\pi$ as the origin, i.e. $(\boldsymbol{\alpha}_\pi = 0, \boldsymbol{\beta}_\pi = 0)$ is not a local minimum for $L(\cdot, \cdot)$. Then there exists a unit direction $\mathbf{d} = (\mathbf{u}, \mathbf{v}) \in \mathbb{R}^{D-d}$ with $\|\mathbf{d}\|^2 = \|\mathbf{u}\|^2 + \|\mathbf{v}\|^2 = 1$ and an $0 < \varepsilon_0 < 1$ such that

$$L(\varepsilon \mathbf{d}) < L(0), \quad \forall 0 < \varepsilon \leq \varepsilon_0 < 1. \quad (13)$$

Clearly $\|\mathbf{u}\| > 0$, otherwise it will contradict Fact 2. Using the definition of directional-derivative, we have that

$$\begin{aligned} \mathbf{d}^\top \nabla^{(2)} L(0, 0) \mathbf{d} &= \lim_{\varepsilon \rightarrow 0} \frac{\langle \nabla L(\varepsilon \mathbf{d}), \mathbf{d} \rangle - \langle \nabla L(0), \mathbf{d} \rangle}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0} \frac{\langle \nabla L(\varepsilon \mathbf{d}), \mathbf{d} \rangle}{\varepsilon}. \end{aligned}$$

On the other hand, using the Hessian structure the LHS equals $\mathbf{d}^\top \nabla^{(2)} L(0, 0) \mathbf{d} \geq \lambda \|\mathbf{u}\|^2 \triangleq K_1 > 0$. Thus

$$\lim_{\varepsilon \rightarrow 0} \frac{\langle \nabla L(\varepsilon \mathbf{d}), \mathbf{d} \rangle}{\varepsilon} = K_1 > 0.$$

Thus there exists an $\varepsilon_1 > 0$ and $K > 0$ such that

$$\frac{\langle \nabla L(\varepsilon \mathbf{d}), \mathbf{d} \rangle}{\varepsilon} \geq K, \quad \forall 0 < \varepsilon \leq \varepsilon_1,$$

which implies

$$\langle \nabla L(\varepsilon \mathbf{d}), \mathbf{d} \rangle \geq K\varepsilon, \quad \forall 0 \leq \varepsilon \leq \varepsilon_1.$$

Defining the function $g : \mathbb{R}_+ \rightarrow \mathbb{R}$ as $g(\varepsilon) = L(\varepsilon \mathbf{d})$, we obtain that $g'(\varepsilon) = \langle \nabla L(\varepsilon \mathbf{d}), \mathbf{d} \rangle \geq K\varepsilon$ for $0 \leq \varepsilon \leq \varepsilon_1$. Using the fundamental theorem of Calculus, we have that for any $0 \leq \varepsilon \leq \varepsilon_1$,

$$\begin{aligned} g(\varepsilon) - g(0) &= \int_0^\varepsilon g'(t) dt \\ &\geq \int_0^\varepsilon Kt dt \\ &= \frac{K\varepsilon^2}{2}. \end{aligned}$$

Thus $g(\varepsilon) = L(\varepsilon \mathbf{d}) \geq L(0) + \frac{K\varepsilon^2}{2}$ for all $0 \leq \varepsilon \leq \varepsilon_1$ whereas $L(\varepsilon \mathbf{d}) < L(0)$ for all $0 < \varepsilon < \varepsilon_0$ from Eq. (13). Choosing $\varepsilon_\star = \min(\varepsilon_0, \varepsilon_1)$, we have a contradiction for $0 < \varepsilon < \varepsilon_\star$. Thus $0 \equiv \boldsymbol{\theta}_\pi = (\boldsymbol{\alpha}_\pi, \boldsymbol{\beta}_\pi)$ is a local minimum. \square

B.5. Proof of Thm. 3

Proof. Since $\bar{\theta}_\pi \triangleq (\theta_\pi, \mathbf{a}_\pi) \in \mathbb{R}^D$ is a canonical extension of $\theta_\pi = (e_\pi = \mathbf{a}_\pi, \dots, b_\pi) \in \mathbb{R}^{D-d}$, which is a local-minimum for $L(\cdot)$ in \mathbb{R}^{D-d} , following the same-steps for the gradient computation and probability evaluation as in proof of Thm. 2 in § B.4, it immediately follows that $\bar{\theta}_\pi$ also satisfies properties (ii)-(iv), i.e. it's a stationary point, it captures the marginal, since $\mathbb{P}_{\bar{\theta}_\pi}(x_{n+1} = 1 | x_1^n) = \mathbb{P}_{\theta_\pi}(x_{n+1} = 1 | x_1^n) = \mathbb{P}(x_{n+1} = 1) = \pi$, and hence its loss equals entropy of stationary distribution $H(\pi)$. In a similar fashion, the Hessian computation is essentially the same except for a slight difference in the Hessian structure, i.e.

$$\mathbf{H}(\bar{\theta}_\pi) \triangleq \nabla^{(2)} L(\bar{\theta}_\pi) = \begin{bmatrix} \mathbf{H}_\alpha & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \mathbf{H}_\alpha = \pi_0 \pi_1 \begin{bmatrix} 1 & \mathbf{u}^\top & 0 \\ \mathbf{u} & \mathbf{V} & (p+q-1)\mathbf{I} \\ 0 & (p+q-1)\mathbf{I} & 0 \end{bmatrix},$$

where $\mathbf{u} \triangleq \sum_{n \in [N]} \frac{\mathbf{p}_n}{N}$, $\mathbf{V} \triangleq \sum_{n \in [N]} \frac{\mathbf{p}_n \mathbf{p}_n^\top}{N}$. In the weight-tied case, we observe that the matrix \mathbf{V} also contains the $(p+q-1)\mathbf{I}$ terms which in the non-weight-tied case gets de-coupled (due to separate e and \mathbf{a} parameters). In fact, \mathbf{H}_α corresponds to the Hessian w.r.t the parameters $\alpha = (b, \mathbf{a}, e)$, i.e. $\mathbf{H}_\alpha = \nabla_\alpha^{(2)} L|_{\alpha=\alpha_\pi}$. Now it remains to show that \mathbf{H}_α is indefinite and hence $\bar{\theta}_\pi$ a saddle point.

Clearly, \mathbf{H}_α cannot be negative definite since with $\mathbf{d} = (1, 0, \dots, 0)$, we have $\mathbf{d}^\top \mathbf{H}_\alpha \mathbf{d} = \pi_0 \pi_1 > 0$ for $(p, q) \in (0, 1)$. Now we show that it cannot be positive definite either. Denoting

$$\mathbf{H}_\alpha = \pi_0 \pi_1 \begin{bmatrix} 1 & \mathbf{b}^\top \\ \mathbf{b} & \mathbf{C} \end{bmatrix}, \quad \mathbf{b} \triangleq \begin{bmatrix} \sum_{n \in [N]} \frac{\mathbf{p}_n}{N} \\ 0 \end{bmatrix}, \quad \mathbf{C} \triangleq \begin{bmatrix} \sum_{n \in [N]} \frac{\mathbf{p}_n \mathbf{p}_n^\top}{N} & (p+q-1)\mathbf{I} \\ (p+q-1)\mathbf{I} & 0 \end{bmatrix}.$$

Using the characterization of positive-definiteness by Schur's complement (Horn and Johnson, 2012), we have that $\mathbf{H}_\alpha \succ 0 \Leftrightarrow 1 > 0$ and $\mathbf{C} - \mathbf{b}\mathbf{b}^\top \succ 0$. This can be further simplified to

$$\begin{aligned} \mathbf{M} \triangleq \mathbf{C} - \mathbf{b}\mathbf{b}^\top &= \begin{bmatrix} \sum_{n \in [N]} \frac{\mathbf{p}_n \mathbf{p}_n^\top}{N} & (p+q-1)\mathbf{I} \\ (p+q-1)\mathbf{I} & 0 \end{bmatrix} - \begin{bmatrix} \sum_{n \in [N]} \frac{\mathbf{p}_n}{N} \\ 0 \end{bmatrix} \begin{bmatrix} \sum_{n \in [N]} \frac{\mathbf{p}_n^\top}{N} & 0 \end{bmatrix} \\ &= \begin{bmatrix} \text{Cov}(\{\mathbf{p}_n\}_{n=1}^N) & (p+q-1)\mathbf{I} \\ (p+q-1)\mathbf{I} & 0 \end{bmatrix}, \end{aligned}$$

where $\text{Cov}(\{\mathbf{p}_n\}_{n=1}^N) = \sum_{n \in [N]} \frac{1}{N} \left(\mathbf{p}_n - \sum_{n \in [N]} \frac{\mathbf{p}_n}{N} \right) \left(\mathbf{p}_n - \sum_{n \in [N]} \frac{\mathbf{p}_n}{N} \right)^\top$ is the covariance matrix of the set $\{\mathbf{p}_n\}_{n=1}^N$. Now we show that \mathbf{M} cannot be positive definite. Suppose not. Then there exists a $\lambda > 0$ such that $\mathbf{v}^\top \mathbf{M} \mathbf{v} \geq \lambda \|\mathbf{v}\|^2$ for all $\mathbf{v} = (\mathbf{v}_1, \mathbf{v}_2) \in \mathbb{R}^{2d}$. This further implies that

$$\mathbf{v}_1^\top \text{Cov}(\{\mathbf{p}_n\}_{n=1}^N) \mathbf{v}_1 + 2(p+q-1) \langle \mathbf{v}_1, \mathbf{v}_2 \rangle \geq \lambda \|\mathbf{v}\|^2, \quad \forall \mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^d.$$

Taking $\mathbf{v}_1 = 0$ the above inequality implies that $\lambda \|\mathbf{v}_2\|^2 \leq 0$ for all $\mathbf{v}_2 \in \mathbb{R}^d$, which is a contradiction. Hence \mathbf{M} cannot be positive definite and consequently neither can \mathbf{H}_α . \square

B.6. Proof of Lemma 1

Proof. Consider the loss function $L(\cdot)$ given in Eq. (1). We can rewrite it as follows:

$$\begin{aligned} L(\bar{\theta}) &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [x_{n+1} \cdot \log f_{\bar{\theta}}(x_1^n) + (1 - x_{n+1}) \cdot \log(1 - f_{\bar{\theta}}(x_1^n))] \\ &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} \left[\mathbb{E}_{x_{n+1}|x_1^n} [x_{n+1}] \cdot \log f_{\bar{\theta}}(x_1^n) + \mathbb{E}_{x_{n+1}|x_1^n} [1 - x_{n+1}] \cdot \log(1 - f_{\bar{\theta}}(x_1^n)) \right] \\ &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} \left[\mathbb{P}(x_{n+1} = 1 | x_n) \log \frac{f_{\bar{\theta}}(x_1^n)}{\mathbb{P}(x_{n+1} = 1 | x_n)} + \mathbb{P}(x_{n+1} = 0 | x_n) \log \frac{1 - f_{\bar{\theta}}(x_1^n)}{\mathbb{P}(x_{n+1} = 0 | x_n)} \right. \\ &\quad \left. - \mathbb{P}(x_{n+1} = 1 | x_n) \log \frac{1}{\mathbb{P}(x_{n+1} = 1 | x_n)} - \mathbb{P}(x_{n+1} = 0 | x_n) \log \frac{1}{\mathbb{P}(x_{n+1} = 0 | x_n)} \right] \end{aligned}$$

$$\begin{aligned}
 &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} \left[\mathbb{P}(x_{n+1} = 1 | x_n) \log \frac{f_{\bar{\theta}}(x_1^n)}{\mathbb{P}(x_{n+1} = 1 | x_n)} + \mathbb{P}(x_{n+1} = 0 | x_n) \log \frac{1 - f_{\bar{\theta}}(x_1^n)}{\mathbb{P}(x_{n+1} = 0 | x_n)} \right] \\
 &\quad - \mathbb{E}_{x_n} \left[\mathbb{E}_{x_1^n | x_n} \left[\mathbb{P}(x_{n+1} = 1 | x_n) \log \frac{1}{\mathbb{P}(x_{n+1} = 1 | x_n)} + \mathbb{P}(x_{n+1} = 0 | x_n) \log \frac{1}{\mathbb{P}(x_{n+1} = 0 | x_n)} \right] \right] \\
 &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} \left[\mathbb{P}(x_{n+1} = 1 | x_n) \log \frac{f_{\bar{\theta}}(x_1^n)}{\mathbb{P}(x_{n+1} = 1 | x_n)} + \mathbb{P}(x_{n+1} = 0 | x_n) \log \frac{1 - f_{\bar{\theta}}(x_1^n)}{\mathbb{P}(x_{n+1} = 0 | x_n)} \right] \\
 &\quad - \mathbb{E}_{x_n} \left[\mathbb{P}(x_{n+1} = 1 | x_n) \log \frac{1}{\mathbb{P}(x_{n+1} = 1 | x_n)} + \mathbb{P}(x_{n+1} = 0 | x_n) \log \frac{1}{\mathbb{P}(x_{n+1} = 0 | x_n)} \right].
 \end{aligned}$$

Since $f_{\bar{\theta}}(x_1^n) = \mathbb{P}_{\bar{\theta}}(x_{n+1} = 1 | x_n)$, we have that the first term above is

$$\begin{aligned}
 &\mathbb{P}(x_{n+1} = 1 | x_n) \log \frac{f_{\bar{\theta}}(x_1^n)}{\mathbb{P}(x_{n+1} = 1 | x_n)} + \mathbb{P}(x_{n+1} = 0 | x_n) \log \frac{1 - f_{\bar{\theta}}(x_1^n)}{\mathbb{P}(x_{n+1} = 0 | x_n)} \\
 &= -D_{\text{KL}}(\mathbb{P}(x_{n+1} = \cdot | x_n) \parallel \mathbb{P}_{\bar{\theta}}(x_{n+1} = \cdot | x_1^n)).
 \end{aligned}$$

Further, observe that the second term is exactly the entropy rate $H(x_{n+1}|x_n)$. Hence, the above expression for the loss reduces to

$$L(\bar{\theta}) = \frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} [D_{\text{KL}}(\mathbb{P}(x_{n+1} = \cdot | x_n) \parallel \mathbb{P}_{\bar{\theta}}(x_{n+1} = \cdot | x_1^n))] + H(x_{n+1}|x_n),$$

and we are done. \square

B.7. Proof of Lemma 2

Proof. It suffices to show that for any component $\bar{\theta}_j$ of $\bar{\theta} \in \mathbb{R}^D$,

$$\begin{aligned}
 \frac{\partial}{\partial \bar{\theta}_j} L(\bar{\theta}) &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} \left[(x_{n+1} - f_{\bar{\theta}}(x_1^n)) \cdot \frac{\partial}{\partial \bar{\theta}_j} (\mathbf{a}^\top \mathbf{z}_n + b) \right] \\
 &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} \left[(\mathbb{P}(x_{n+1} = 1 | x_n) - f_{\bar{\theta}}(x_1^n)) \cdot \frac{\partial}{\partial \bar{\theta}_j} (\mathbf{a}^\top \mathbf{z}_n + b) \right].
 \end{aligned}$$

Recall from Eq. (3) that $L(\cdot)$ is given by

$$L(\bar{\theta}) = -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} [x_{n+1} \cdot \log f_{\bar{\theta}}(x_1^n) + (1 - x_{n+1}) \cdot \log(1 - f_{\bar{\theta}}(x_1^n))],$$

which implies that

$$\begin{aligned}
 \frac{\partial}{\partial \bar{\theta}_j} L(\bar{\theta}) &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} \left[x_{n+1} \cdot \frac{\partial}{\partial \bar{\theta}_j} \log f_{\bar{\theta}}(x_1^n) + (1 - x_{n+1}) \cdot \frac{\partial}{\partial \bar{\theta}_j} \log(1 - f_{\bar{\theta}}(x_1^n)) \right] \\
 &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} \left[x_{n+1} \cdot \frac{1}{f_{\bar{\theta}}(x_1^n)} \frac{\partial}{\partial \bar{\theta}_j} f_{\bar{\theta}}(x_1^n) + (1 - x_{n+1}) \cdot \frac{1}{1 - f_{\bar{\theta}}(x_1^n)} \frac{\partial}{\partial \bar{\theta}_j} (1 - f_{\bar{\theta}}(x_1^n)) \right].
 \end{aligned}$$

Since $f_{\bar{\theta}}(x_1^n) = \sigma(\mathbf{a}^\top \mathbf{z}_n + b)$, we first note that the derivative of σ is given by $\sigma'(z) = \frac{e^{-z}}{(1+e^{-z})^2} = \sigma(z)(1 - \sigma(z))$. Hence, the derivative $\frac{\partial}{\partial \bar{\theta}_j} f_{\bar{\theta}}(x_1^n)$ can be written as

$$\begin{aligned}
 \frac{\partial}{\partial \bar{\theta}_j} f_{\bar{\theta}}(x_1^n) &= \frac{\partial}{\partial \bar{\theta}_j} \sigma(\mathbf{a}^\top \mathbf{z}_n + b) = \sigma(\mathbf{a}^\top \mathbf{z}_n + b) [1 - \sigma(\mathbf{a}^\top \mathbf{z}_n + b)] \frac{\partial}{\partial \bar{\theta}_j} (\mathbf{a}^\top \mathbf{z}_n + b) \\
 &= f_{\bar{\theta}}(x_1^n) [1 - f_{\bar{\theta}}(x_1^n)] \frac{\partial}{\partial \bar{\theta}_j} (\mathbf{a}^\top \mathbf{z}_n + b).
 \end{aligned}$$

Plugging this into the above expression, we have

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} L(\bar{\theta}) &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} \left[x_{n+1} \cdot [1 - f_{\bar{\theta}}(x_1^n)] \frac{\partial}{\partial \theta_j} (\mathbf{a}^\top \mathbf{z}_n + b) - (1 - x_{n+1}) \cdot f_{\bar{\theta}}(x_1^n) \frac{\partial}{\partial \theta_j} (\mathbf{a}^\top \mathbf{z}_n + b) \right] \\
 &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} \left[x_{n+1} \cdot [1 - f_{\bar{\theta}}(x_1^n)] \frac{\partial}{\partial \theta_j} (\mathbf{a}^\top \mathbf{z}_n + b) - (1 - x_{n+1}) \cdot f_{\bar{\theta}}(x_1^n) \frac{\partial}{\partial \theta_j} (\mathbf{a}^\top \mathbf{z}_n + b) \right] \\
 &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} \left[(x_{n+1} - f_{\bar{\theta}}(x_1^n)) \cdot \frac{\partial}{\partial \theta_j} (\mathbf{a}^\top \mathbf{z}_n + b) \right] \\
 &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} \left[(\mathbb{E}_{x_{n+1}|x_1^n} [x_{n+1}] - f_{\bar{\theta}}(x_1^n)) \cdot \frac{\partial}{\partial \theta_j} (\mathbf{a}^\top \mathbf{z}_n + b) \right] \\
 &= -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^n} \left[(\mathbb{P}(x_{n+1} = 1 | x_n) - f_{\bar{\theta}}(x_1^n)) \cdot \frac{\partial}{\partial \theta_j} (\mathbf{a}^\top \mathbf{z}_n + b) \right],
 \end{aligned}$$

and we are done. □

C. Additional results for first-order Markov chains

C.1. Model architecture and hyper-parameters

Table 1: Parameters in the transformer architecture with their shape.

| Parameter | Matrix shape |
|---|---------------|
| transformer.wte | $2 \times d$ |
| transformer.wpe | $N \times d$ |
| transformer.h.ln_1 ($\times \ell$) | $d \times 1$ |
| transformer.h.attn.c_attn ($\times \ell$) | $3d \times d$ |
| transformer.h.attn.c_proj ($\times \ell$) | $d \times d$ |
| transformer.h.ln_2 ($\times \ell$) | $d \times 1$ |
| transformer.h.mlp.c_fc ($\times \ell$) | $4d \times d$ |
| transformer.h.mlp.c_proj ($\times \ell$) | $d \times 4d$ |
| transformer.ln_f | $d \times 1$ |

Table 2: Settings and parameters for the transformer model used in the experiments.

| | |
|---------------------|---|
| Dataset | k -th order binary Markov source |
| Architecture | Based on the GPT-2 architecture as implemented in (Pagliardini, 2023) |
| Batch size | Grid-searched in $\{16, 50\}$ |
| Accumulation steps | 1 |
| Optimizer | AdamW ($\beta_1 = 0.9, \beta_2 = 0.95$) |
| Learning rate | 0.001 |
| Scheduler | Cosine |
| # Iterations | 8000 |
| Weight decay | 1×10^{-3} |
| Dropout | 0 |
| Sequence length | Grid-searched in $\{512, 1024, 2048\}$ |
| Embedding dimension | Grid-searched in $\{4, 8, 16, 32, 64\}$ |
| Transformer layers | Between 1 and 6 depending on the experiment |
| Attention heads | Grid-searched in $\{1, 2, 4, 8\}$ |
| Mask window | Between 2 and full causal masking depending on the experiment |
| Repetitions | 3 or 5 |

D. Empirical formula for $p + q < 1$ based on low-rank solutions

In this section we compute the function $f_{\theta}(x_1^n)$ that gives the next-symbol probability predicted by the network, using the values of the weight matrices obtained five independent experiment runs. By substituting the empirical weights into the transformer architecture from § A, i.e.

$$\begin{aligned}
 \mathbf{x}_n &= x_n \mathbf{e} + \mathbf{p}_n \in \mathbb{R}^d, && \text{(Uni-embedding)} \\
 \mathbf{y}_n &= \mathbf{x}_n + \mathbf{W}_O \sum_{i \in [n]} \text{att}_{n,i} \cdot \mathbf{W}_V \mathbf{x}_i \in \mathbb{R}^d, && \text{(Attention)} \\
 \mathbf{z}_n &= \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) \in \mathbb{R}^d, && \text{(FF)} \\
 \text{logit}_n &= \langle \mathbf{a}, \mathbf{z}_n \rangle + b \in \mathbb{R}, && \text{(Linear)} \\
 f_{\theta}(x_1^n) &\triangleq \mathbb{P}_{\theta}(x_{n+1} = 1 \mid x_1^n) = \sigma(\text{logit}_n). && \text{(Prediction)}
 \end{aligned}$$

We can obtain an explicit expression for $f_{\theta}(x_1^n)$ as it is actually learned by the model. We now analyze each section of the model architecture separately.

Embedding. All the five independent runs show that the word embedding vector e has the structure

$$e = e \cdot v \quad (14)$$

where $v = (v_1, \dots, v_d)$ is such that $v_i \in \{-1, +1\}$ for all i , i.e., $v \in \{-1, 1\}^d$, and e is some constant. Moreover, the positional embeddings are approximately constant across positions n , and they share a similar structure to e . In particular, we always have that

$$p_n = p = p \cdot v \quad \forall n \quad (15)$$

for some constant $p \in \mathbb{R}$. Furthermore, the constants are always such that $p < 0$ and $e + p > 0$.

Attention. Across all the runs, we observe that the contribution of the attention mechanism is negligible compared to the skip-connection. In particular, we observe that

$$\frac{\|\mathbf{W}_O \sum_{i \in [n]} \text{att}_{n,i} \cdot \mathbf{W}_V x_i\|}{\|\mathbf{y}_n\|} \approx 0.01 \quad (16)$$

uniformly for all n . Therefore, we can use the approximation

$$\mathbf{y}_n \approx \mathbf{x}_n \quad \forall n. \quad (17)$$

FF. For the MLP layer, we observe that \mathbf{W}_1 and \mathbf{W}_2 have a clear joint structure. In fact, we empirically see that

$$\mathbf{W}_1 = w_1 \cdot \mathbf{w} \cdot \mathbf{v}^T \quad (18)$$

where \mathbf{v} is again the same vector as in Eq. (14), $\mathbf{w} \in \{-1, 1\}^r$ and $w_1 \in \mathbb{R}$. Hence, \mathbf{W}_1 is a rank-one matrix. As customary in the GPT-2 model, for our experiments we used $r = 4d = 16$. Furthermore, we see that

$$\mathbf{W}_2 = \mathbf{W}_1^T. \quad (19)$$

Due to this structure and the formula for \mathbf{y}_n described above, we have

$$\mathbf{W}_1 \mathbf{y}_n = \mathbf{W}_1 \mathbf{x}_n = w_1 d(e x_n + p) \mathbf{w} \quad (20)$$

Let now $\mathbf{r} = \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n)$. Due to the fact that $p < 0$ and $e + p > 0$, we have that, if $x_n = 1$,

$$r_i = \begin{cases} e + p, & \text{if } w_i = 1, \\ 0, & \text{if } w_i = -1. \end{cases} \quad (21)$$

While if $x_n = 0$,

$$r_i = \begin{cases} 0, & \text{if } w_i = 1, \\ -p, & \text{if } w_i = -1. \end{cases} \quad (22)$$

Let $\beta = \sum_{i=1}^r \mathbb{1}_{\{w_i=1\}}$. Since $\mathbf{W}_2 = \mathbf{W}_1^T = w_1 \mathbf{v} \cdot \mathbf{w}^T$, we have that, for $\tilde{\mathbf{r}} = \mathbf{W}_2 \mathbf{r}$,

$$\tilde{\mathbf{r}} = \begin{cases} w_1^2 d(e + p) \beta \cdot \mathbf{v}, & \text{if } x_n = 1, \\ w_1^2 d p (r - \beta) \cdot \mathbf{v}, & \text{if } x_n = 0. \end{cases} \quad (23)$$

Or more compactly,

$$\tilde{\mathbf{r}} = w_1^2 d(e x_n + p) ((2\beta - r)x_n + r - \beta) \cdot \mathbf{v}, \quad (24)$$

and

$$\mathbf{z}_n = \mathbf{y}_n + \tilde{\mathbf{r}} = (e x_n + p) (1 + w_1^2 d ((2\beta - r)x_n + r - \beta)) \cdot \mathbf{v} \quad (25)$$

Linear. Since $\mathbf{a} = e$ due to weight-tying, we have

$$\text{logit}_n = e d(e x_n + p) (1 + w_1^2 d ((2\beta - r)x_n + r - \beta)) + b \quad (26)$$

Prediction. We can now plug in the empirical values obtained by averaging five independent runs. The numerical results that we get are

$$\begin{aligned}
 e &= 0.3618 \\
 p &= -0.1539 \\
 w_1 &= 0.3264 \\
 b &= -0.1229 \\
 \beta &= 5
 \end{aligned}
 \tag{27}$$

Plugging these numbers into Eq. (26), we get

$$\text{logit}_n = \begin{cases} 0.8191, & \text{if } x_n = 1, \\ -1.3897, & \text{if } x_n = 0. \end{cases}
 \tag{28}$$

Hence, by applying the sigmoid function to the logit values, we obtain the predicted probabilities

$$f_{\theta}(x_n) = \mathbb{P}_{\theta}(x_{n+1} = 1 \mid x_n) = \begin{cases} \sigma(0.8191) = 0.694, & \text{if } x_n = 1, \\ \sigma(-1.3897) = 0.199, & \text{if } x_n = 0. \end{cases}
 \tag{29}$$

The numerical results correspond almost exactly to the expected theoretical values of $1 - q = 0.7$ and $p = 0.2$.