
Local to Global: Learning Dynamics and Effect of Initialization for Transformers

Anonymous Authors¹

Abstract

In recent years, transformer-based models have revolutionized deep learning, particularly in sequence modeling. To better understand this phenomenon, there is a growing interest in using Markov input processes to study transformers. However, our current understanding in this regard remains limited with many fundamental questions about how transformers learn Markov chains still unanswered. In this paper, we address this by focusing on first-order Markov chains and single-layer transformers, providing a comprehensive characterization of the learning dynamics in this context. Specifically, we prove that transformer parameters trained on next-token prediction loss can either converge to global or local minima, contingent on the initialization and the Markovian data properties, and we characterize the precise conditions under which this occurs. To the best of our knowledge, this is the first result of its kind highlighting the role of initialization. We further demonstrate that our theoretical findings are corroborated by empirical evidence. Based on these insights, we provide guidelines for the initialization of transformer parameters and demonstrate their effectiveness. Finally, we outline several open problems in this arena. Code is available at: <https://anonymous.4open.science/r/Local-to-Global-C70B/>.

1. Introduction

Transformers have been at the forefront of recent successes across various fields including natural language processing (Vaswani et al., 2017). To obtain insights into their impressive sequential modeling capabilities, a notable emerging theme among several recent works is to model

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

the input data as a Markov process.

Using this Markovian perspective, works such as (Nichani et al., 2024; Edelman et al., 2024; Bietti et al., 2023), among others, study the in-context learning capabilities of a two-layer transformer. (Makkuva et al., 2024) analyzes the loss-landscape for the next-token prediction task, while (Ildiz et al., 2024) shows an equivalence between the attention mechanism and Markov models. Although these works reveal interesting insights about transformers and their capabilities, many fundamental questions about their learning dynamics remain unanswered. In particular, a comprehensive characterization of their training dynamics vis-à-vis the data distributional properties and the role of initialization is still missing.

To address this gap, in this paper, we focus on the canonical setting of first-order Markov chains and single-layer transformers and make the following contributions:

- **Theoretical analysis:** We precisely characterize the loss landscape and gradient flow dynamics for single-layer transformers with first-order Markov chains (Secs. 3 and 4). We demonstrate that transformer parameters trained on next-token prediction loss can converge to global or local minima, depending on the initialization and the Markovian data properties, and determine the exact conditions under which this occurs (Thms. 2, 3, and 8). To the best of our knowledge, this is the first result of its kind.
- **Insights into initialization:** Our theoretical analysis underscores the crucial role of initialization in transformer parameter training. Specifically, we demonstrate how the standard Gaussian initialization scheme can lead the convergence to local or global minima depending on the Markovian data properties (Thms. 2 and 8, Figs. 1 and 2).
- **Guidelines:** Based on these insights, we provide practical guidelines for parameter initialization, corroborated by empirical evidence demonstrating their effectiveness (§ A.3.2).

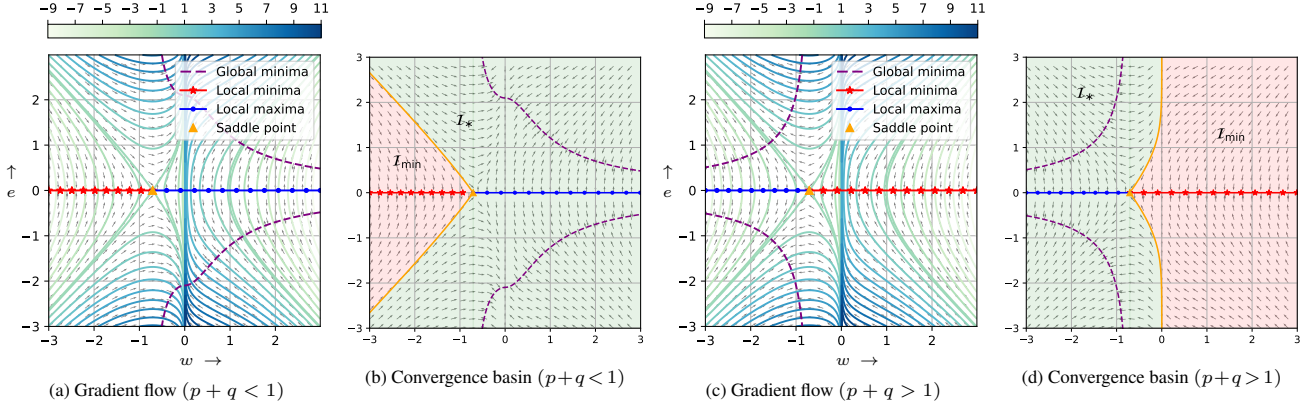


Figure 1: Gradient flow dynamics and initialization effect for single-layer transformers. (p, q) are Markov switching probabilities, and (e, w) are the embedding and weight parameters (Sec. 2). (a), (c): The flow is aligned along energy contour lines, converging to local or global optima. (b), (d): \mathcal{I}_* is the basin of convergence for global minima, \mathcal{I}_{\min} for the local minima, and yellow asymptotes for the saddle point. Notice the contrasting behavior for Gaussian initialization around origin for $p+q \leq 1$.

2. Problem setting

We succinctly define the problem setting for analysis of single-layer transformers with Markovian data (§ A).

Input data. We assume that the input word sequence $\{x_n\}_{n=1}^N \sim (\pi, \mathbf{P})$ is a first-order time-homogenous Markov chain with binary state space $\{0, 1\}$ and a fixed kernel $\mathbf{P} = [1-p, p; q, 1-q]$, where $p = \mathbf{P}_{01} = \mathbb{P}(x_{n+1} = 1 | x_n = 0)$ and $q = \mathbf{P}_{10} = \mathbb{P}(x_{n+1} = 0 | x_n = 1)$ denote the switching probabilities from the states 0 and 1 respectively. We call $p+q$ the *switching factor*. $\pi \triangleq (\pi_0, \pi_1) = (q, p)/(p+q)$ is the stationary distribution satisfying $\pi = \pi \mathbf{P}$.

Transformer architecture. We consider a single-layer transformer with a single-head attention and ReLU non-linearity, which for an input sequence $\{x_n\}_{n=1}^N$, performs the following mathematical operations at each $n \in [N]$:

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

$$\mathbf{y}_n = \mathbf{x}_n + \sum_{i \in [n]} \underbrace{\text{att}_{n,i}}_{\in (0,1)} \cdot \mathbf{W}_V \mathbf{x}_i \in \mathbb{R}^d, \quad (\text{Attention})$$

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

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

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

Here $\theta \triangleq (\mathbf{e}, \{\mathbf{p}_n\}_{n=1}^N, \dots, \mathbf{W}_1, \mathbf{W}_2, b, \mathbf{a}) \in \mathbb{R}^D$ denotes the full list of the transformer parameters (§ B).

Loss and training. The transformer parameters θ are trained using gradient-based methods to minimize the cross-

entropy loss on the next-token prediction, L , given by

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

Loss landscape. A key surprising observation in (Makkuva et al., 2024) is that the loss function $L(\cdot)$ admits both global and local minima depending on the switching factor $p+q$ of the Markovian data, and the weight-tying of the embedding and linear weights ($\mathbf{e} = \mathbf{a}$) of the transformer. In view of these results, we focus on the weight-tying scenario and hence let $\mathbf{e} = \mathbf{a}$ to be a single parameter in \mathbb{R}^d . Thus, $\theta = (\mathbf{e} = \mathbf{a}, \{\mathbf{p}_n\}_{n=1}^N, \dots, \mathbf{W}_1, \mathbf{W}_2, b)$. We interchangeably refer to θ as both the transformer and the set of parameters.

Our objective. While the aforementioned results detail the static landscape of the loss, they do not characterize the learning dynamics on the loss surface and the effect of initialization, which plays a central role in training machine learning models (Arora et al., 2019). In view of these shortcomings, the main objective of this paper is to address the following question:

(Q.1): Can we explain how the initialization and learning dynamics affect the convergence of the transformer parameters to the local or global optima?

3. Canonical low-rank parameterization

Motivation. Given the complexity of the transformer architecture and the non-convex loss function, it is challenging to analyze the learning dynamics directly (Edelman et al., 2024). To tackle this, we capitalize on the following empiri-

cal observation (§ A.3.1) which is the motivating idea behind our approach: when trained by gradient-based methods, the weight matrices ($\mathbf{W}_V, \dots, \mathbf{W}_1, \mathbf{W}_2$) at the optima θ_* and θ_{\min} exhibit *rank-one* structure, whose eigenvector is the same direction in which the both the token embedding e and the positional embeddings p_n are all aligned in. Interestingly, such low-rank solutions can also be shown to be theoretically optimal (see § C). While these observations illustrate the implicit bias towards low-rank solutions at the final convergence, a natural question arises: *if we initialize with low-rank parameters, will they remain low-rank during training?* In § A.3.1, we affirmatively address this based on a thorough empirical evaluation for single-layer transformers and inspired by these empirical phenomena, without loss of generality, we restrict our attention to these low-rank manifolds to characterize the learning dynamics.

Parameterization. More specifically, we consider a special low-rank parameterization that is empirically observed and capitalize on it to address (Q.1). Interestingly, along this low-rank manifold, it suffices to consider a reduced set of parameters $\theta \in \mathbb{R}^2$ or $\theta \in \mathbb{R}^3$ given by:

$$\theta = (e, w) \in \mathbb{R}^2, \text{ or } \theta = (e, w, a) \in \mathbb{R}^3. \quad (2)$$

Here e denotes the *embedding* scalar, w the *weight*, and a the *attention* parameter respectively. These parameters distill the essential role of the embedding vectors ($e, \{p_n\}$) in the [Embedding](#) layer, the weight matrices ($\mathbf{W}_1, \mathbf{W}_2$) in the [FF](#) layer, and the attention matrices $\mathbf{W}_{K,Q,V}$ in the [Attention](#) layer respectively. We refer to § A.1 for a more detailed description. For the ease of exposition, here we let $a = 0$ and analyze the learning dynamics for $\theta = (e, w) \in \mathbb{R}^2$ and defer the general case $\theta \in \mathbb{R}^3$ to § A.2.

Using this parameterization $\theta = (e, w)$ in the transformer architecture (Sec. 2) and the equivalence between the cross-entropy loss and the logistic loss $\ell_{\log}(z) = \log(1 + \exp(-z))$, $z \in \mathbb{R}$, the loss function in Eq. (1) can be compactly written as (Lemma 6):

$$L(\theta) = \frac{1}{N} \sum_{n \in [N]} \mathbb{E}[\ell_{\log}((2x_{n+1} - 1) \cdot \text{logit}_n(\theta))], \quad (3)$$

$$\text{logit}_n(\theta) = \left(e^2(1 + 2w|w|)x_n + b_* - \frac{e^2}{2} \right),$$

where $b_*(\theta) = \arg\min_{b \in \mathbb{R}} L(\theta, b)$ is the optimal bias minimizing the loss for each θ and has a closed form expression (Lemma 5). Empirically, this roughly translates to running the gradient-descent for the bias for more steps at each θ . In practice, one additional step is usually sufficient (§ A.3).

3.1. Loss landscape with canonical parameterization

With the new set of parameters $\theta = (e, w)$, we are now ready to analyze the loss $L(\cdot)$ in Eq. (3). First we recall the

definition of a critical point (Lee et al., 2016). A point θ_* is a critical point for L if $\nabla L(\theta_*) = 0$. Similarly, the standard notions of local & global minima and saddle points (§ A.1). Equipped with these definitions, Thm. 1 below characterizes the loss landscape in terms of these local and global optima.

Theorem 1 (All critical points). *Let the input sequence be $\{x_n\}_{n=1}^N \sim (\pi, \mathbf{P})$, the transformer parameters $\theta = (e, w)$, and the next-token prediction loss L be as in Eq. (3). Then for any $(p, q) \in (0, 1)^2$ with $p + q \neq 1$ and $N \in \mathbb{N}$,*

(i) $\Theta_*(p, q)$, the set of all global minima is given by

$$\{(e, w) : e^2(1 + 2w|w|) = \log(1 - p)(1 - q)/pq\},$$

(ii) $\Theta_{\min}(p, q)$, the set of all local minima is given by

$$\{(e, w) : e = 0, (p + q - 1)(1 + 2w|w|) > 0\},$$

(iii) $\Theta_{\max}(p, q)$, the set of all local maxima is given by

$$\{(e, w) : e = 0, (p + q - 1)(1 + 2w|w|) < 0\},$$

(iv) and $\Theta_{\text{sad}}(p, q)$, the set of all saddle points is

$$\{(0, -1/\sqrt{2})\}.$$

Thus the set of all critical points is

$$\{\theta : \nabla L(\theta) = 0\} = \Theta_* \cup \Theta_{\min} \cup \Theta_{\max} \cup \Theta_{\text{sad}}.$$

In addition, for any $\theta_* \in \Theta_*$, $\theta_{\min} \in \Theta_{\min}$, $\theta_{\max} \in \Theta_{\max}$, and $\theta_{\text{sad}} \in \Theta_{\text{sad}}$, the loss values satisfy

$$H(x_{n+1} | x_n) = L(\theta_*) < L(\theta_{\min}) = L(\theta_{\max}) \\ = L(\theta_{\text{sad}}) = H(x_{n+1}).$$

Proof. We refer to § F. □

Fig. 1 illustrates the loci of these critical points for $p + q \leq 1$. Motivated by empirical observations, while (Makkuva et al., 2024) characterizes local minima for $p + q > 1$, it is interesting to note that our Thm. 1 shows that local minima also exist for $p + q < 1$ (Fig. 1a). So why did they observe the minima only for the former? The answer to this, and more broadly to question (Q.1) lies in the learning dynamics for θ , influenced by initialization, which we study in the next section.

4. Learning dynamics

Capitalizing on the loss landscape in Thm. 1, we now focus on the convergence of gradient-based algorithms to these critical points. Specifically, we focus on the gradient-flow of the parameters, $(\theta_t)_{t \geq 0}$, governed by

$$\frac{d\theta_t}{dt} = -\nabla L(\theta_t), \quad \theta_t = (e_t, w_t) \in \mathbb{R}^2, t \geq 0, \quad (\text{GF})$$

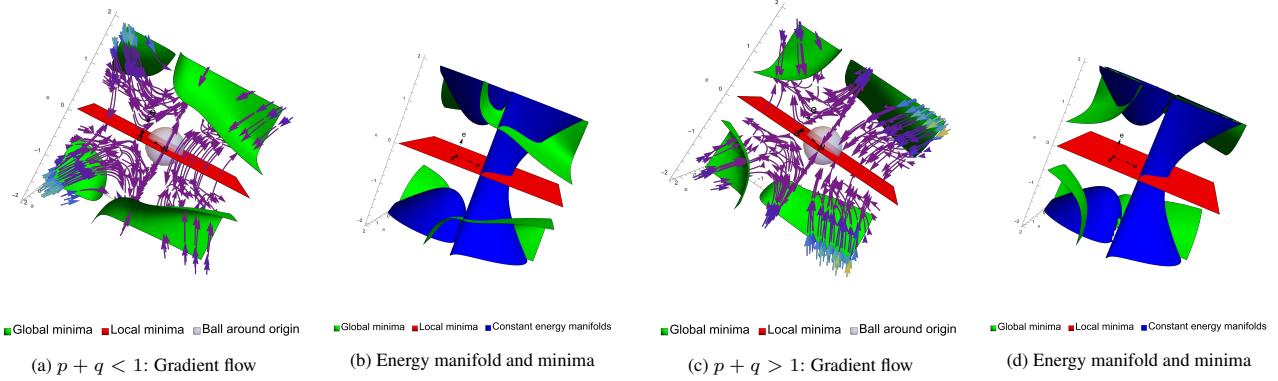


Figure 2: Gradient flow dynamics for the canonical parameters $\theta = (e, w, a) \in \mathbb{R}^3$ with the attention scalar a (§ A.2). Notice the contrasting behavior for Gaussian initialization around origin for $p + q$ smaller and greater than one. For an enhanced view of the flow near the origin, please refer to Fig. 5.

where $\theta_t \triangleq \theta(t)$ is a continuously differentiable curve in \mathbb{R}^2 starting with a randomly initialized θ_0 . To characterize these trajectories, we define an *energy function* $\mathcal{E}(\cdot, \cdot)$, which plays a crucial role in the GF dynamics. It is defined as

$$\mathcal{E}(e, w) \triangleq e^2 - (w^2 + \text{sign}(w) \cdot \log |w|), \quad (4)$$

for all $(e, w) \in \mathbb{R}^2 \setminus \text{e-axis}$, where $\text{e-axis} \triangleq \{(e, w = 0)\}$. Figs. 1a and 1c illustrate these energy contour lines. The utility of the energy function is captured below.

Lemma 1 (Constant energy along the flow). *For any $(p, q) \in (0, 1)^2$ and initialization $\theta_0 = (e_0, w_0)$, let $(\theta_t)_{t \geq 0}$ be the corresponding GF trajectory starting from θ_0 . If $\theta_0 \in \mathbb{R}^2 \setminus \text{e-axis}$, the energy stays constant along the trajectory, i.e. for all $t \geq 0$,*

$$\mathcal{E}(\theta_t) = e_t^2 - (w_t^2 + \text{sign}(w_t) \cdot \log |w_t|) = \mathcal{E}(\theta_0). \quad (5)$$

On the other hand, if $\theta_0 \in \text{e-axis}$, we have $\theta_t \in \text{e-axis}$ for all $t \geq 0$ with $w_t = w_0 = 0$.

We are now ready to present the main results of our paper. Specifically, Thm. 2 and Thm. 8, for $p + q \geq 1$ respectively, highlight the role of the switching factor and the parameter initialization, θ_0 , in deciding whether the GF converges to local or global optima. First we define the energy $\mathcal{E}_{\text{sad}} \triangleq \mathcal{E}(e = 0, w = -1/\sqrt{2}) = -(1 + \log 2)/2$.

Theorem 2 (GF dynamics for $p + q > 1$). *Let $(p, q) \in (0, 1)^2$ with $p + q > 1$, the input sequence be $\{x_n\}_{n=1}^N \sim (\pi, \mathbf{P})$, and $(\theta_t)_{t \geq 0}$ be the corresponding GF trajectory starting from θ_0 . Then for all initializations $\theta_0 \in \mathbb{R}^2$, the gradient flow converges to a critical point of the loss L . That is, there exists a $\theta_{\text{lim}} \in \mathbb{R}^2$ such that $\lim_{t \rightarrow \infty} \theta_t = \theta_{\text{lim}}$ and $\nabla L(\theta_{\text{lim}}) = 0$. In particular, θ_{lim} is a*

(i) *a local minimum if*

$$\begin{aligned} \theta_0 \in \mathcal{I}_{\text{min}} \triangleq \{ & (e, w) : w \in (-1/\sqrt{2}, 0), \\ & e \in (-g(w), g(w)), g(w) = \sqrt{w^2 - \log(-w) + \mathcal{E}_{\text{sad}}} \\ & \cup \{(e, w) : w \geq 0\}, \end{aligned}$$

- (ii) *a saddle point if $\theta_0 \in \mathcal{I}_{\text{sad}} \triangleq \{(e, w) : w \in [-1/\sqrt{2}, 0), e = \pm \sqrt{w^2 - \log(-w) + \mathcal{E}_{\text{sad}}}\}$,*
 (iii) *a local maximum if $\theta_0 \in \mathcal{I}_{\text{max}} \triangleq \{(e, w) : e = 0, w < -1/\sqrt{2}\}$,*
 (iv) *and a global minimum if $\theta_0 \in \mathcal{I}_* \triangleq \mathbb{R}^2 \setminus (\mathcal{I}_{\text{min}} \cup \mathcal{I}_{\text{sad}} \cup \mathcal{I}_{\text{max}})$.*

Consequently, when $p + q > 1$, if we use the standard initialization $\theta_0 \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_2)$ with $\sigma^2 \ll 1/\sqrt{2}$, θ_{lim} will be a local minimum with high probability. If $p + q < 1$, under the same initialization scheme, θ_{lim} will be a global minimum with high probability.

Key insights and conclusion. Together, Thms. 2 and 8 address our motivating question (Q.1) by fully characterizing the GF dynamics in terms of initialization and input data properties. Specifically, our results explain the phenomenon in (Makkuva et al., 2024) wherein they observe local minima for $p + q > 1$ more often than for $p + q < 1$, owing to standard Gaussian initialization around origin (Figs. 1b and 1d). However, in practice, we often do not know the input switching factor, raising a natural question: *is there a data-agnostic initialization that always converges to global minima?* Indeed, as can be seen from Figs. 1b and 1d, there is a common region of initialization, $\mathcal{I}_{\text{common}} \triangleq \{(e, w) : w < 0, |e| > \sqrt{w^2 - \log(-w) + \mathcal{E}_{\text{sad}}}\}$, that leads to the global minima convergence irrespective of the switching (§ A.3.2). We believe our findings open interesting avenues of future research for GF analysis with deeper architectures and higher order Markov chains.

References

- 220
221
222 Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Su-
223 vrit Sra. Transformers learn to implement preconditioned
224 gradient descent for in-context learning. In *Thirty-seventh*
225 *Conference on Neural Information Processing Systems*,
226 2023. URL [https://openreview.net/forum?](https://openreview.net/forum?id=LziniAXEI9)
227 [id=LziniAXEI9](https://openreview.net/forum?id=LziniAXEI9).
- 228 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu
229 Ma, and Denny Zhou. What learning algorithm is in-
230 context learning? Investigations with linear models. In
231 *The Eleventh International Conference on Learning Rep-*
232 *resentations*, 2023.
- 233
234 Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei
235 Hu. A convergence analysis of gradient descent for deep
236 linear neural networks, 2019.
- 237
238 Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song
239 Mei. Transformers as statisticians: Provable in-context
240 learning with in-context algorithm selection. In *Workshop*
241 *on Efficient Systems for Foundation Models @ ICML2023*,
242 2023.
- 243
244 Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve
245 Jegou, and Leon Bottou. Birth of a transformer: A mem-
246 ory viewpoint. In *Thirty-seventh Conference on Neural*
247 *Information Processing Systems*, 2023.
- 248
249 Pritam Chandra, Tanmay Kumar Sinha, Kabir Ahuja, Ankit
250 Garg, and Navin Goyal. Towards analyzing self-attention
251 via linear neural network, 2024. URL [https://](https://openreview.net/forum?id=4fVuBf5HE9)
252 openreview.net/forum?id=4fVuBf5HE9.
- 253
254 Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran
255 Yang. Training dynamics of multi-head softmax attention
256 for in-context learning: Emergence, convergence, and
257 optimality, 2024.
- 258
259 Thomas M Cover and Joy A Thomas. *Elements of informa-*
260 *tion theory*. John Wiley & Sons, 2nd edition, 2006.
- 261
262 John M Danskin. The theory of max-min, with applications.
263 *SIAM Journal on Applied Mathematics*, 14(4):641–664,
264 1966.
- 265
266 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong
267 Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and
268 Zhifang Sui. A survey on in-context learning, 2023. URL
269 <https://arxiv.org/abs/2301.00234>.
- 270
271 Benjamin L. Edelman, Ezra Edelman, Surbhi Goel, Eran
272 Malach, and Nikolaos Tsilivis. The evolution of statisti-
273 cal induction heads: In-context learning markov chains,
274 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.](https://transformer-circuits.pub/2021/framework/index.html)
[pub/2021/framework/index.html](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.
- M Emrullah Ildiz, Yixiao Huang, Yingcong Li, Ankit Singh
Rawat, and Samet Oymak. From self-attention to markov
models: Unveiling the dynamics of generative transform-
ers. *arXiv preprint arXiv:2402.13512*, 2024.
- Samy Jelassi, Michael Eli Sander, and Yuanzhi Li. Vision
transformers provably learn spatial structure. In Alice H.
Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun
Cho, editors, *Advances in Neural Information Processing*
Systems, 2022.
- Jason D Lee, Max Simchowitz, Michael I Jordan, and Ben-
jamin Recht. Gradient descent only converges to minimiz-
ers. In *Conference on learning theory*, pages 1246–1257.
PMLR, 2016.
- Yingcong Li, M. Emrullah Ildiz, Dimitris Papailiopoulos,
and Samet Oymak. Transformers as algorithms: general-
ization and stability in in-context learning. In *Proce-*
edings of the 40th International Conference on Machine
Learning, 2023a.
- Yuchen Li, Yuanzhi Li, and Andrej Risteski. How do trans-
formers learn topic structure: Towards a mechanistic
understanding, 2023b.
- Ashok Vardhan Makkuva, Marco Bondaschi, Adway Girish,
Alliot Nagle, Martin Jaggi, Hyeji Kim, and Michael Gast-
par. Attention with Markov: A framework for princi-
pled analysis of transformers via Markov chains. *arXiv*
preprint arXiv:2402.04161, 2024.
- Eshaan Nichani, Alex Damian, and Jason D Lee. How
transformers learn causal structure with gradient descent.
arXiv preprint arXiv:2402.14735, 2024.

- 275 Samet Oymak, Ankit Singh Rawat, Mahdi Soltanolkotabi,
276 and Christos Thrampoulidis. On the role of attention in
277 prompt-tuning. In *Proceedings of the 40th International*
278 *Conference on Machine Learning*, 2023.
- 279 Matteo Pagliardini. GPT-2 modular codebase implementa-
280 tion. <https://github.com/epfml/llm-baselines>, 2023.
- 282 Jorge Pérez, Pablo Barceló, and Javier Marinkovic. Attention
283 is Turing-complete. *Journal of Machine Learning*
284 *Research*, 22(75):1–35, 2021.
- 285 Charlie Snell, Ruiqi Zhong, Dan Klein, and Jacob Steinhardt.
286 Approximating how single head attention learns. *CoRR*,
287 abs/2103.07601, 2021. URL <https://arxiv.org/abs/2103.07601>.
- 289 Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya
290 Gunasekar, and Nathan Srebro. The implicit bias of
291 gradient descent on separable data. *Journal of Machine*
292 *Learning Research*, 19(70):1–57, 2018.
- 293 Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang,
294 and Samet Oymak. Max-margin token selection in atten-
295 tion mechanism. In *Thirty-seventh Conference on Neural*
296 *Information Processing Systems*, 2023.
- 299 Yuandong Tian, Yiping Wang, Beidi Chen, and Si-
300 mon Shaolei Du. Scan and snap: Understanding training
301 dynamics and token composition in 1-layer transformer.
302 In *Conference on Parsimony and Learning (Recent Spot-*
303 *light Track)*, 2023.
- 305 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszko-
306 reit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and
307 Illia Polosukhin. Attention is all you need. In *Advances*
308 *in Neural Information Processing Systems*, pages 5998–
309 6008, 2017.
- 310 Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo,
311 João Sacramento, Alexander Mordvintsev, Andrey Zh-
312 moginov, and Max Vladymyrov. Transformers learn in-
313 context by gradient descent. In *International Conference*
314 *on Machine Learning*, pages 35151–35174, 2023.
- 316 Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo,
317 João Sacramento, Alexander Mordvintsev, Andrey Zh-
318 moginov, and Max Vladymyrov. Transformers learn in-
319 context by gradient descent, 2023.
- 320 Colin Wei, Yining Chen, and Tengyu Ma. Statistically
321 meaningful approximation: a case study on approxim-
322 ating Turing machines with transformers. In *Advances*
323 *in Neural Information Processing Systems*, volume 35,
324 pages 12071–12083, 2022.
- 326 Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like
327 transformers. In *International Conference on Machine*
328 *Learning*, pages 11080–11090, 2021.
- 329 Sang Michael Xie, Aditi Raghunathan, Percy Liang,
and Tengyu Ma. An explanation of in-context learn-
ing 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 uni-
versal approximators of sequence-to-sequence functions?
In *International Conference on Learning Representations*,
2020.

Contents

330	1 Introduction	1
331		
332	2 Problem setting	2
333		
334	3 Canonical low-rank parameterization	2
335	3.1 Loss landscape with canonical parameterization	3
336		
337	4 Learning dynamics	3
338		
339	A Supporting results to the manuscript	8
340	A.1 Canonical low-rank parameterization	9
341	A.2 Gradient flow with attention	10
342	A.3 Empirical results	10
343	A.3.1 Low-rank parameters	11
344	A.3.2 Effect of initialization: broader implications	12
345	A.4 Related works	12
346		
347	B Single-layer transformer: architecture and results	14
348	B.1 Loss landscape results	14
349		
350	C Low-rank structure of the optima	15
351		
352	D Canonical reparameterization	16
353		
354	E Analysis of the loss with the bias, $L(\theta, b)$, in Eq. (8) and Eq. (26)	18
355	E.1 Technical lemmas	18
356		
357	F Analysis of the loss without bias, $L(\theta)$, and proof of Thm. 1	20
358	F.1 Proof of Thm. 1	21
359		
360	G Proofs of Thm. 2 and Thm. 8	22
361	G.1 Proof of Thm. 2	22
362	G.2 GF dynamics for $p + q < 1$	24
363		
364	H Gradient flow analysis with attention	26
365	H.1 Canonical parameterization with attention	26
366	H.2 Analysis of the loss function $L(\theta)$ from Eq. (45)	28
367	H.3 Gradient flow analysis	30
368	H.4 Role of Standard Initialization	30
369		
370		
371		
372		
373		
374		
375		
376		
377		
378		
379		
380		
381		
382		
383		
384		

385	I Proofs of theorems in App. E	31
386	I.1 Proof of Thm. 7	31
387		
388	I.2 Proof of Thm. 6	33
389		
390	J Proofs of technical lemmas in App. E	34
391		
392	J.1 Proof of Lemma 2	34
393	J.2 Proof of Lemma 3	35
394		
395	J.3 Proof of Lemma 4	36
396		
397	K Proofs of lemmas in App. F	39
398		
399	K.1 Proof of Lemma 5	39
400	K.2 Proof of Lemma 6	39
401		
402	K.3 Proof of Lemma 7	39
403	K.4 Proof of Lemma 8	39
404		
405	L Proofs of lemmas in App. G	40
406		
407	L.1 Proof of Lemma 9	40
408	L.2 Proof of Lemma 10	40
409		
410	L.3 Proof of Lemma 11	42
411		
412	L.4 Proof of Lemma 12	42
413		
414	M Proofs of lemmas in App. H	43
415		
416	M.1 Proof of Lemma 13	43
417	M.2 Optimality conditions for linear self- attention	45
418	M.3 Stationary points	46
419		
420	M.3.1 Stationary Points Where $\mathbb{E}_X [f_1 X + f_2] = 0, e = 0$	48
421	M.3.2 Stationary Points Where $\mathbb{E}_X [f_1 X + f_2] = 0, e \neq 0, 1 + ae^2 = 0, 1 + 2w w = 0$	51
422	M.3.3 Computing the Optimal Bias	51
423		
424	M.4 Proof of 15	53
425		
426	M.5 Role of Standard Initialization	54
427		
428	N Additional empirical results	57
429		
430		
431	O Model architecture and hyper-parameters	58
432		

A. Supporting results to the manuscript

We provide additional theoretical and empirical details for our main results in the manuscript. First we define the notation.

Notation. We denote scalars by italic lower case letters like x, y and Euclidean vectors and matrices in bold: $\mathbf{x}, \mathbf{y}, \mathbf{M}$, etc. $\|\cdot\|$ denotes 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 $\sigma(z) \triangleq 1/(1 + e^{-z})$,

ReLU(z) \triangleq $\max(0, z)$ and the convex logistic loss $\ell_{\log}(z) \triangleq \log(1 + \exp(-z)) \in (0, \infty)$. 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 $H(x) = H(\mathbf{p}_x) \triangleq -\sum_{i \in [k]} p_i \log p_i$. The conditional entropy is defined to be $H(y|x) \triangleq H(x, y) - H(x)$. The entropy rate of a stochastic process $(x_n)_{n \geq 1}$ is defined as $\lim_{n \rightarrow \infty} H(x_1^n)/n$. We simply write $x = y$ to mean $\mathbb{P}(x = y) = 1$. We also use the shorthand $\mathbb{P}(y = j | x)$ for $\mathbb{P}(y = j | x = x)$ as a function of the random variable x . For $p \in (0, 1)$, the binary entropy function $h(\cdot)$ is defined as $h(p) \triangleq -p \log p - (1 - p) \log(1 - p)$.

A.1. Canonical low-rank parameterization

Recall from Sec. 3 that the reduced set of parameters $\boldsymbol{\theta} \in \mathbb{R}^2$ or $\boldsymbol{\theta} \in \mathbb{R}^3$ is given by:

$$\boldsymbol{\theta} = (e, w) \in \mathbb{R}^2, \text{ or } \boldsymbol{\theta} = (e, w, a) \in \mathbb{R}^3. \quad (\text{Reparameterization})$$

Here e denotes the *embedding* scalar, w the *weight*, and a the *attention* parameter respectively. Now we describe the parameterization of the transformer vis-à-vis these scalars and refer to App. D for a more detailed description. Let the input $\{x_n\}_{n=1}^N$ be a first-order Markov chain as in Sec. 2 and let $n \in [N]$ be fixed. Then we have

$$\begin{aligned} \text{Embedding} : e = e \cdot \boldsymbol{\alpha}, \mathbf{p}_n = \left(-\frac{e}{2}\right) \cdot \boldsymbol{\alpha} \rightarrow \mathbf{x}_n = e \left(x_n - \frac{1}{2}\right) \boldsymbol{\alpha}, \quad e \in \mathbb{R}, \boldsymbol{\alpha} \in \{\pm 1\}^d / \sqrt{d}, \\ \text{Attention} : \mathbf{W}_V = \boldsymbol{\alpha} \mathbf{v}^\top \rightarrow \mathbf{y}_n = e \left(x_n - \frac{1}{2}\right) \boldsymbol{\alpha} + \underbrace{\langle \mathbf{v}, \boldsymbol{\alpha} \rangle}_{\triangleq a \approx 0} \left(\sum_{i \in [n]} \text{att}_{n,i} \cdot e \left(x_i - \frac{1}{2}\right)\right) \boldsymbol{\alpha}, \mathbf{v} \in \mathbb{R}^d. \end{aligned}$$

The scalar a is the product of $\langle \mathbf{v}, \boldsymbol{\alpha} \rangle$ and the scaling in the attention weights $\text{att}_{n,i}$, which is empirically close to zero for first-order Markov chains. Hence for the ease of exposition, we first omit it by letting $a = 0$ and analyze the general case when $a \in \mathbb{R}$ in App. A.2. We continue:

$$\text{FF} : \mathbf{W}_1 = \frac{|w|}{\sqrt{d}} \mathbf{1} \boldsymbol{\alpha}^\top, \mathbf{W}_2 = \frac{w}{\sqrt{d}} \boldsymbol{\alpha} \mathbf{1}^\top \rightarrow \mathbf{z}_n = e \left(x_n - \frac{1}{2}\right) (1 + 4w|w|x_n) \boldsymbol{\alpha}, w \in \mathbb{R}.$$

$\mathbf{1}$ is the all-one vector in \mathbb{R}^r with $r = 4d$ typically in practice. Substituting this \mathbf{z}_n in the linear layer with $e = a$ and bias $b \in \mathbb{R}$, the logits and the probabilities simplify to:

$$\text{Linear} : \text{logit}_n(e, w, b) = e^2(1 + 2w|w|x_n) + b - \frac{e^2}{2} \in \mathbb{R}, \quad (6)$$

$$\text{Prediction} : f_{(\boldsymbol{\theta}, b)}(x_1^n) = \sigma(\text{logit}_n) \in (0, 1), \quad \boldsymbol{\theta} \triangleq (e, w). \quad (7)$$

Finally, using the equivalence between the cross-entropy loss and the logistic loss $\ell_{\log}(\cdot)$, the loss function in Eq. (1) can be compactly written as (Lemma 6):

$$L(\boldsymbol{\theta}, b) = \frac{1}{N} \sum_{n \in [N]} \mathbb{E}[\ell_{\log}((2x_{n+1} - 1) \cdot \text{logit}_n(\boldsymbol{\theta}))], \quad \boldsymbol{\theta} \in \mathbb{R}^2, b \in \mathbb{R}. \quad (8)$$

Due to convexity of $\ell_{\log}(\cdot)$, it follows that $L(\boldsymbol{\theta}, b)$ is convex in the bias b for any fixed $\boldsymbol{\theta}$, whose minimizer, $b_\star(\boldsymbol{\theta}) = \text{argmin}_{b \in \mathbb{R}} L(\boldsymbol{\theta}, b)$, has a closed form expression (Lemma 5). Hence, without loss of generality, we consider the loss with this optimal bias b_\star :

$$L(\boldsymbol{\theta}) \triangleq L(\boldsymbol{\theta}, b_\star) = \frac{1}{N} \sum_{n \in [N]} \mathbb{E} \left[\ell_{\log} \left((2x_{n+1} - 1) \left(e^2(1 + 2w|w|x_n) + b_\star - \frac{e^2}{2} \right) \right) \right]. \quad (9)$$

Empirically, this roughly translates to running the gradient-based algorithm for the bias for more steps at each $\boldsymbol{\theta}$. In practice, one additional step is usually sufficient (see App. A.3). Eq. (9) resembles the standard logistic regression loss (Soudry et al., 2018) whose binary labels are $2x_{n+1} - 1 \in \{\pm 1\}$ and the logits given by $e^2(1 + 2w|w|x_n) + b_\star - e^2/2$, for each $n \in [N]$. The key difference here is that the logits are a non-linear function of the parameters (e, w) unlike in the standard setting.

Definitions of local and global optima. For the sake of completeness, we first recall the definition of a critical point (Lee et al., 2016). A point $\theta_* \in \mathbb{R}^2$ is a critical or a stationary point for L if $\nabla L(\theta_*) = 0$. A critical point θ_* is a *local minimum* if there exists a neighborhood U around θ_* such that $L(\theta_*) \leq L(\theta)$ for all $\theta \in U$, and a *local maximum* if $L(\theta_*) \geq L(\theta)$. If the neighborhood U is whole of \mathbb{R}^2 , it is a *global minimum/maximum*. On the other hand, a critical point is a saddle point if for all neighborhoods U around θ_* , there are $\theta_1, \theta_2 \in U$ such that $L(\theta_1) \leq L(\theta_*) \leq L(\theta_2)$.

A.2. Gradient flow with attention

In this section, we consider the attention scalar $a \in \mathbb{R}$ (Sec. 3) and study the gradient flow dynamics with the parameters $\theta = (e, w, a) \in \mathbb{R}^3$. The parameter a captures the overall scaling from the value, key, and query components in the attention layer. Recall that the soft-max attention weights are given by $\text{att}_{n,i} \propto \exp(\langle \mathbf{q}_n, \mathbf{k}_i \rangle / \sqrt{d})$, where $\mathbf{q}_n = \mathbf{W}_Q \mathbf{x}_n$ and $\mathbf{k}_i = \mathbf{W}_K \mathbf{x}_i$ are the query and key embeddings for any position $i \in [n]$. Using the low-rank structure of the query and key matrices, satisfying $\mathbf{W}_Q^\top \mathbf{W}_K = (q^2 d) \alpha \alpha^\top$ and the value matrix $\mathbf{W}_V = \alpha \mathbf{v}^\top$ for some $q \in \mathbb{R}$ and $\mathbf{v} \in \mathbb{R}^d$ (App. H), and assuming linear attention $\text{att}_{n,i} \propto \langle \mathbf{q}_n, \mathbf{k}_i \rangle / \sqrt{d}$, we define a single scalar $a \triangleq \langle \mathbf{v}, \alpha \rangle q^2 d^{5/2} / 4$ that captures the essence of the attention layer. We note that linear attention weights are a standard assumption in the transformer analysis literature (Ahn et al., 2023; von Oswald et al., 2023). Using this parameterization, similar to the steps in Sec. 3, we obtain the final loss function to be

$$L(\theta) = \mathbb{E} \left[\ell_{\log} \left((2Y - 1) \left(e^2 \left[\left(X - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2)| \right] + b_* \right) \right) \right],$$

where $\theta = (e, w, a)$ and b_* is the corresponding optimal bias. L recovers the loss in Eq. (9) when $a = 0$. In Thm. 10, we determine the set of all critical points of L in terms of global minima and local optima in closed-form expressions, analogous to Thm. 1. Capitalizing on this characterization, we now shift our focus to the analysis of the gradient flow in \mathbb{R}^3 . To this end, let $(\theta_t)_{t \geq 0}$ be a C^1 (continuously differentiable) curve in \mathbb{R}^3 governed by

$$\frac{d\theta_t}{dt} = -\nabla L(\theta_t), \quad \theta_t = (e_t, w_t, a_t) \in \mathbb{R}^3, \quad t \geq 0, \quad (\text{GF-attn})$$

starting with a randomly initialized θ_0 . We define the *energy function* $\mathcal{E}(\cdot, \cdot, \cdot)$ as

$$\mathcal{E}(e, w, a) \triangleq e^2 - (w^2 + \text{sign}(w) \cdot \log |w|) - 2a^2, \quad \forall (e, w, a) \in \mathbb{R}^3 \setminus \text{ea-plane}, \quad (10)$$

where ea-plane $\triangleq \{(e, w = 0, a)\}$. It is similar to its counterpart in Eq. (4), except for the $2a^2$ term. Fig. 2 visualizes this energy surface and the set of critical points, which reveal close resemblance to that of Fig. 1 in \mathbb{R}^2 . Capitalizing on the energy function, we now present our main result with the attention.

Theorem 3 (GF dynamics with attention). *For any $(p, q) \in (0, 1)^2$ and initialization $\theta_0 \in \mathbb{R}^3$, let $(\theta_t)_{t \geq 0}$ be the corresponding GF-attn trajectory starting from it. Then for all $\theta_0 \in \mathbb{R}^3$, the gradient flow converges to a critical point of the loss L . That is, there exists a $\theta_{\text{lim}} \in \mathbb{R}^3$ such that $\lim_{t \rightarrow \infty} \theta_t = \theta_{\text{lim}}$ and $\nabla L(\theta_{\text{lim}}) = 0$. Further,*

- (i) *if $\theta_0 \in \mathbb{R}^3 \setminus \text{ea-plane}$, we have $\mathcal{E}(\theta_{\text{lim}}) = \mathcal{E}(\theta_t) = \mathcal{E}(\theta_0)$ for all $t \geq 0$. Hence θ_{lim} is at the intersection of the energy contour line $\mathcal{E} = \mathcal{E}_0$ with that of the set of critical points.*
- (ii) *if $\theta_0 \in \text{ea-plane}$, we have $\theta_t \in \text{ea-plane}$ for all $t \geq 0$ and hence $\theta_{\text{lim}} \in \text{ea-plane}$.*

Thm. 3 shows that the learning dynamics with attention closely resemble those without it (Thms. 2 and 8). While the set of all critical points of L , and thus the limit points of the flow, has a closed-form expression (Thm. 10), deriving the same for the initialization sets \mathcal{I}_{\min} and \mathcal{I}_* to determine the basin of convergence is technically challenging (see discussion in App. H). Nonetheless, empirical observations with the standard Gaussian initialization around origin reveal a similar picture as in the two-dimensional setting for both the $p + q < 1$ and $p + q > 1$ cases (Fig. 2). We believe it's an interesting direction of future research to theoretically characterize this, analogous to Thms. 2 and 8. We refer to App. H for additional details and proofs.

A.3. Empirical results

We empirically validate the low-rank assumption behind our canonical parameterization and investigate our findings about local optima and initialization in the context of the full model.

A.3.1. LOW-RANK PARAMETERS

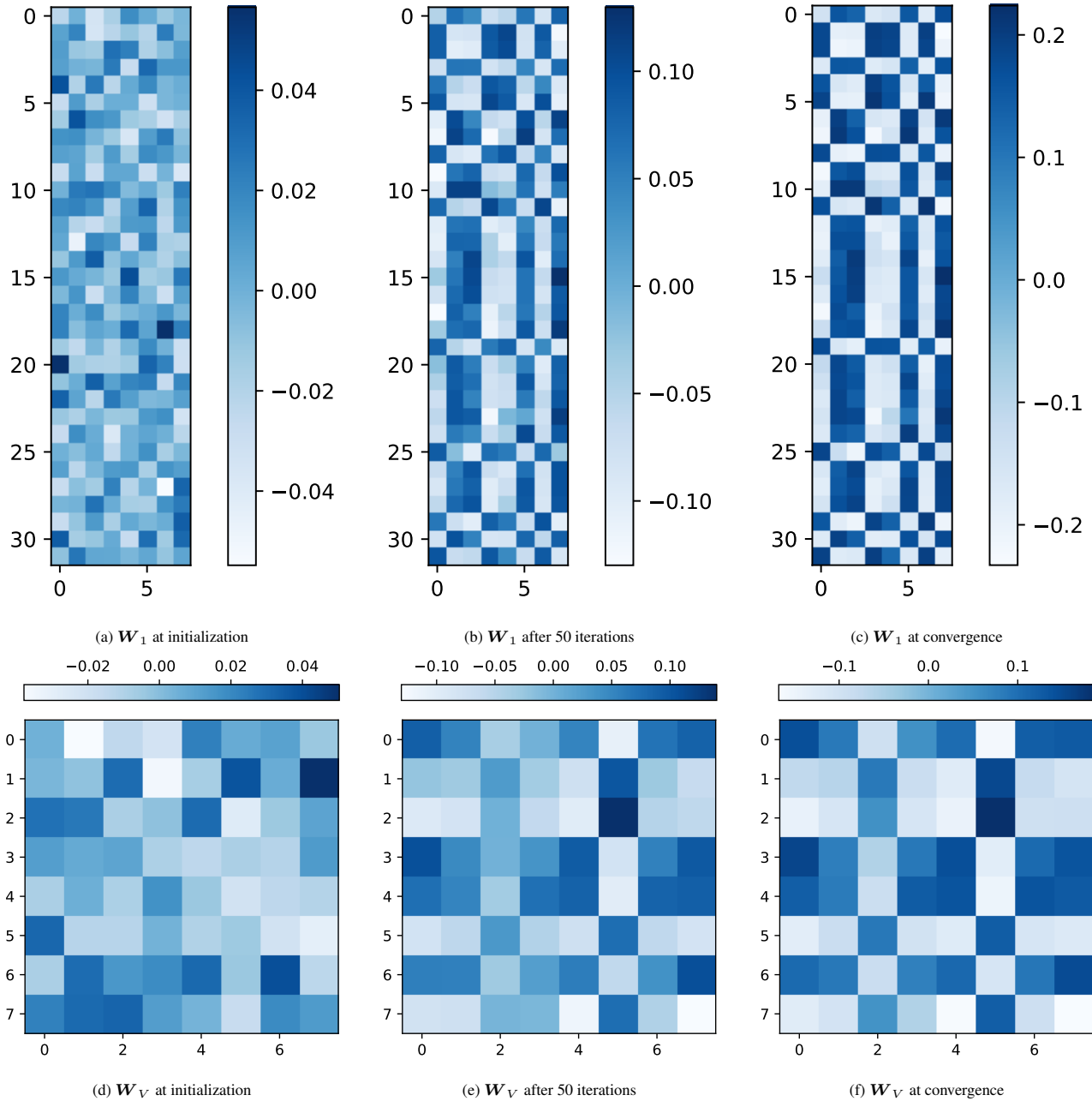


Figure 3: Evolution of parameters W_1 and W_V across iterations, starting from a standard gaussian initialization. At convergence, all the parameter matrices are approximately rank-one.

Low-rank at convergence. We let the input Markov sequence to be $\{x_n\}_{n=1}^N \sim (\pi(p, q), \mathbf{P}(p, q))$ for $p = 0.2, q = 0.3, N = 1024$ and consider the single-layer transformer as defined in Sec. 2 with embedding dimension $d = 8$. First, we initialize the parameters $\theta = (e = \mathbf{a}, \{\mathbf{p}_n\}_{n=1}^N, \dots, \mathbf{W}_1, \mathbf{W}_2, b)$ using the standard Gaussian initialization with standard deviation 0.001 (Pagliardini, 2023) and train them using SGD on a batch size $B = 16$ and for $t = 800$ iterations. In Fig. 3, we track the value matrix $W_V \in \mathbb{R}^{d \times d}$ and the weight matrix $W_1 \in \mathbb{R}^{4d \times d}$ across iterations. We observe that at convergence both W_V and W_1 are approximately rank-one with one of their components being same as the embedding vector (the row in W_V and column in W_1). Further, the embedding vector has all entries in $\{\pm 1\}$ up to a scaling. We observe the same conclusion for other weight matrices $W_{K,Q}, W_2$ and for all values of $(p, q) \in (0, 1)^2$.

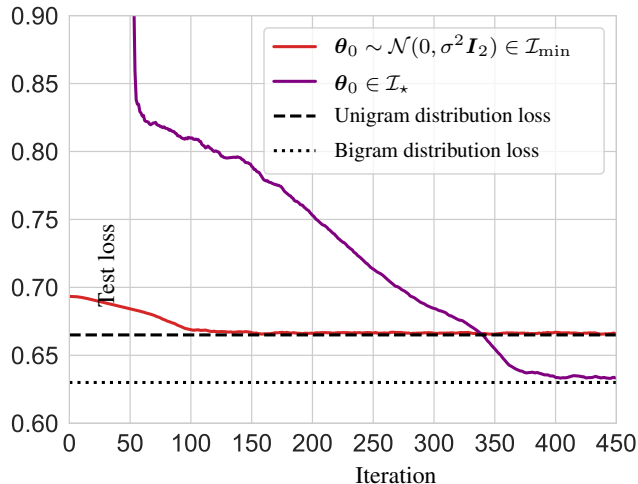


Figure 4: Comparison between the average loss curve for the standard gaussian initialization around 0 and our initialization, for $p = 0.5$ and $q = 0.8$. Starting from the standard initialization, the model converges to a local minimum corresponding to the unigram model. With our initialization, it converges to the global minimum corresponding to the bigram model.

Low-rank initialization remains low-rank during training. Inspired by the low-rank structure obtained above, we randomly initialize the weight parameters as rank-one matrices and the embeddings on the hypercube $\{\pm 1\}^d$. After the initialization, we train them without any low-rank restrictions, and track them during the course of training. Interestingly, here we observe that the parameters still stay low-rank as illustrated in Fig. 6. A similar conclusion holds for the remaining weight matrices. Together these results provide the empirical basis for our canonical parameterization analysis in Sec. 3.

A.3.2. EFFECT OF INITIALIZATION: BROADER IMPLICATIONS

Now we investigate the findings of Sec. 3 and Sec. 4, derived for the canonical low-rank model, more broadly in the context of the general single-layer transformer in Sec. 2. In particular, as shown in Thm. 2 and Fig. 1d for $p + q > 1$, any small initialization around zero would lead a local minima convergence. To test this hypothesis, we compare the standard initialization where all the transformer parameters $\theta = (e = \mathbf{a}, \{\mathbf{p}_n\}_{n=1}^N, \dots, \mathbf{W}_1, \mathbf{W}_2, b)$ are randomly chosen around zero with small variance $\sigma = 0.02$, with a new initialization based on our results, where we initialize the embedding vector e such that all coordinates are equal to $e = 0.5$, \mathbf{W}_1 to be constant with the scalar $w_1 = 1$ and \mathbf{W}_2 constant with $w_2 = -1$ (corresponding to \mathcal{I}_* in Fig. 1d). We indeed observe that the final test loss matches the unigram loss for the standard initialization, while it converges to the optimal bigram loss for our initialization (see Fig. 4). Together these results indicate that though our analysis used canonical parameterization, the corresponding insights are more general and apply more broadly to the general architecture. In a similar spirit, analysis of initialization effects for deeper architectures is an interesting avenue of future research.

A.4. Related works

The recent success of transformer models in deep learning has sparked significant interest and active research in understanding them (Weiss et al., 2021; Oymak et al., 2023; Fu et al., 2023; Pérez et al., 2021; Elhage et al., 2021; Wei et al., 2022; Yun et al., 2020; Tarzanagh et al., 2023). In relation to our paper, they can be broadly classified into two topics: (i) **In-context learning (ICL)**: ICL refers to the ability of transformers learn and reason from information present in their context (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., 2023a; Garg et al., 2022). Along this thread, the works most relevant to ours are (Bietti et al., 2023; Edelman et al., 2024; Nichani et al., 2024), which use Markovian input data to understand the ICL mechanism. (Bietti et al., 2023; Edelman et al., 2024) heuristically show how gradient-based updates can learn an induction-head mechanism using a simplified transformer architecture with frozen encodings, query matrices and linear activations. On the other hand, we consider the transformer in full generality including ReLU nonlinearity, capitalizing on inherent low-rank parameters, to provide a full characterization of the learning dynamics. (Nichani et al., 2024) demonstrates how two-layer transformers with GD learn induction head mechanism when the input has a causal tree dependency, such as in Markov chains. In

660 this work, we focus on the GF dynamics for single-layer transformers and show how they can also converge to local
661 optima, further highlighting the role of initialization. (ii) **Training dynamics:** On the other hand, numerous works have
662 investigated the training dynamics of transformers. For instance, (Chandra et al., 2024) examines the gradient flow in a
663 simplified single-layer transformer, while (Tian et al., 2023) studies the process by which self-attention integrates input
664 tokens, assuming the decoder learns faster than the attention layer. Unlike these settings, our focus is on understanding the
665 training dynamics of the full transformer model without any simplifications. Other related works include (Snell et al., 2021),
666 which analyzes gradient dynamics in LSTM Seq2seq models, (Jelassi et al., 2022), which shows how Vision Transformers
667 learn spatial structures, and (Li et al., 2023b), which demonstrates that a single-layer transformer can learn a constrained
668 topic model. A closely related work is (Ildiz et al., 2024), which shows that self-attention has a Markovian structure, but
669 our focus is on self-attention’s capability in modeling Markov chains and the associated training dynamics.

670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714

B. Single-layer transformer: architecture and results

We first describe the transformer architecture from Sec. 2:

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

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

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

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

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

Here $\boldsymbol{\theta} \triangleq (\mathbf{e}, \{\mathbf{p}_n\}_{n=1}^N, \dots, \mathbf{W}_1, \mathbf{W}_2, b, \mathbf{a}) \in \mathbb{R}^D$ denotes the full list of the transformer parameters from the embedding layer till the linear layer. In the attention layer, 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, $\text{att}_{n,i} \triangleq \text{softmax}(\langle \mathbf{q}_n, \mathbf{k}_i \rangle, \dots, \langle \mathbf{q}_n, \mathbf{k}_n \rangle) / \sqrt{d}$. $\mathbf{W}_{K,Q,V} \in \mathbb{R}^{d \times d}$ are the respective key, query, and value matrices. For multi-headed attention, the same operation is performed on multiple parallel heads, whose outputs are additively combined.

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

$$L(\boldsymbol{\theta}) \triangleq -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [x_{n+1} \cdot \log f_{\boldsymbol{\theta}}(x_1^n) + (1 - x_{n+1}) \cdot \log(1 - f_{\boldsymbol{\theta}}(x_1^n))]. \quad (11)$$

In this paper, we focus on the weight-tied scenario where $\mathbf{e} = \mathbf{a}$. Hence we let them be a single parameter with $\boldsymbol{\theta} = (\mathbf{e} = \mathbf{a}, \{\mathbf{p}_n\}_{n=1}^N, \dots, b) \in \mathbb{R}^D$, where D is the total parameter dimensionality.

B.1. Loss landscape results

Now we recall the theoretical results from (Makkuva et al., 2024) about the loss landscape of L in the form of global and local minima.

Theorem 4 (Global minimum). *Let the input sequence be $\{x_n\}_{n=1}^N \sim (\boldsymbol{\pi}(p, q), \mathbf{P}(p, q))$ for some fixed $(p, q) \in (0, 1)^2$. Then for all (p, q) , there exists a $\boldsymbol{\theta}_* \in \mathbb{R}^D$ with an explicit construction such that it is a global minimum for the population loss $L(\cdot)$ in Eq. (11), i.e.*

$$(i) \quad L(\boldsymbol{\theta}) \geq L(\boldsymbol{\theta}_*) \text{ for all } \boldsymbol{\theta} \in \mathbb{R}^D.$$

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

$$(ii) \quad \mathbb{P}_{\boldsymbol{\theta}_*}(x_{n+1} = 1 \mid x_1^n) = \mathbb{P}(x_{n+1} = 1 \mid x_n), \text{ the Markov kernel.}$$

$$(iii) \quad L(\boldsymbol{\theta}_*) = H(x_{n+1} \mid x_n), \text{ the entropy rate of the Markov chain.}$$

$$(iv) \quad \nabla L(\boldsymbol{\theta}_*) = 0, \text{ i.e. } \boldsymbol{\theta}_* \text{ is a stationary point.}$$

Let $L_* \triangleq L(\boldsymbol{\theta}_*)$ be the global minimal loss from Thm. 4. Now we recall the result on the bad local minimum.

Theorem 5 (Bad local minimum). *Let the input sequence be $\{x_n\}_{n=1}^N \sim (\boldsymbol{\pi}(p, q), \mathbf{P}(p, q))$ for some fixed $(p, q) \in (0, 1)^2$. If $p + q > 1$, there exists an explicit $\boldsymbol{\theta}_{\min} \in \mathbb{R}^D$ such that it is a bad local minimum for the loss $L(\cdot)$, i.e.*

$$(i) \quad \text{there exists a neighborhood } \mathcal{B}(\boldsymbol{\theta}_{\min}, r) \text{ with } r > 0 \text{ such that } L(\boldsymbol{\theta}) \geq L(\boldsymbol{\theta}_{\min}) \text{ for all } \boldsymbol{\theta} \in \mathcal{B}(\boldsymbol{\theta}_{\min}, r), \text{ with } L(\boldsymbol{\theta}_{\min}) > L_*.$$

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

$$(ii) \quad \mathbb{P}_{\boldsymbol{\theta}_{\min}}(x_{n+1} = 1 \mid x_1^n) = \mathbb{P}(x_{n+1} = 1) = \pi_1, \text{ the marginal distribution.}$$

$$(iii) \quad L(\boldsymbol{\theta}_{\min}) = H(x_{n+1}) = H(\boldsymbol{\pi}), \text{ the entropy of the marginal.}$$

$$(iv) \quad \nabla L(\boldsymbol{\theta}_{\min}) = 0, \text{ i.e. } \boldsymbol{\theta}_{\min} \text{ is a stationary point.}$$

C. Low-rank structure of the optima

Here we recall the low-rank structure for the global minima found by SGD consistently across multiple runs when $p + q < 1$ (Appendix C.2) [Makkuva2024attention]. In particular, it is observed that the token and positional encodings point in the same direction α , which is a low-rank factor for the weight matrices in the attention and the feedforward layers, which in turn are all rank-one. Mathematically,

Embedding. The embedding vector e obeys

$$e = e \cdot \alpha$$

for some $e > 0$ and $\alpha \in \{\pm 1\}^d$. Further, the positional embeddings p_n are constant across positions n pointing in the same direction albeit with a negative scalar, i.e.

$$p_n = -p \cdot \alpha$$

for $p > 0$ and $p \approx \frac{e}{2}$ such that $e > p$. Thus from [Embedding](#) layer,

$$x_n = (ex_n - p) \cdot \alpha, \quad (12)$$

which ensures that the respective embeddings for the bit $x_n = 0$ and $x_n = 1$ are $x_n = -p \cdot \alpha$ and $x_n = (e - p) \cdot \alpha$, which are roughly anti-podal.

Attention. Recall from the [Attention](#) layer that the output y_n is given by $y_n = x_n + \mathbf{W}_O \sum_{i \in [n]} \text{att}_{n,i} \cdot \mathbf{W}_V x_i$, where the attention weights $\text{att}_{n,i}$ are computed according to $\text{att}_{n,i} = \exp(\langle \mathbf{q}_n, \mathbf{k}_i \rangle / \sqrt{d}) / \left(\sum_{j \in [n]} \exp(\langle \mathbf{q}_n, \mathbf{k}_j \rangle / \sqrt{d}) \right)$ with $\mathbf{q}_n = \mathbf{W}_Q x_n$ and $\mathbf{k}_i = \mathbf{W}_K x_i$. Here it is observed that the matrix products are all rank-one with α being a factor, i.e.

$$\begin{aligned} \mathbf{W}_O \mathbf{W}_V &= \alpha \cdot \mathbf{v}^\top \in \mathbb{R}^{d \times d}, \quad \text{for some } \mathbf{v} \in \mathbb{R}^d, \\ \mathbf{W}_Q^\top \mathbf{W}_K &= (q^2 d) \alpha \cdot \alpha^\top \in \mathbb{R}^{d \times d}, \quad \text{for some } q \in \mathbb{R}. \end{aligned}$$

Hence,

$$\mathbf{W}_V x_i = \langle \mathbf{v}, \alpha \rangle (ex_i - p) \alpha,$$

and

$$\begin{aligned} \frac{\langle \mathbf{q}_n, \mathbf{k}_i \rangle}{\sqrt{d}} &= \frac{1}{\sqrt{d}} \cdot x_n^\top \mathbf{W}_Q^\top \mathbf{W}_K x_n = \frac{q^2 d}{\sqrt{d}} \cdot (x_n^\top \alpha) (x_i^\top \alpha) \stackrel{(\|\alpha\|^2=d)}{=} \frac{q^2 d^3}{\sqrt{d}} \cdot (ex_n - p)(ex_i - p) \\ &= q^2 d^{5/2} \cdot (ex_n - p)(ex_i - p). \end{aligned}$$

Thus,

$$\begin{aligned} y_n &= x_n + \sum_{i \in [n]} \text{att}_{n,i} \cdot \mathbf{W}_O \mathbf{W}_V x_i \\ &= (ex_n - p) \alpha + \sum_{i \in [n]} \text{att}_{n,i} \cdot \langle \mathbf{v}, \alpha \rangle (ex_i - p) \alpha \\ &= \left((ex_n - p) + \langle \mathbf{v}, \alpha \rangle \sum_{i \in [n]} \frac{\exp(q^2 d^{5/2} (ex_n - p)(ex_i - p))}{\sum_{j \in [n]} \exp(q^2 d^{5/2} (ex_n - p)(ex_j - p))} \cdot (ex_i - p) \right) \alpha. \end{aligned} \quad (13)$$

It is further noticed that $\langle \mathbf{v}, \alpha \rangle \approx 0$ and hence $y_n = (ex_n - p) \alpha = x_n$.

Feed-forward. For the [Feed-forward](#) layer, both the matrices $\mathbf{W}_1 \in \mathbb{R}^{r \times d}$ and $\mathbf{W}_2 \in \mathbb{R}^{d \times r}$ exhibit rank-one structure with α being one of the factors,

$$\mathbf{W}_1 = w \cdot \mathbf{w} \cdot \alpha^\top, \quad \text{for some } \mathbf{w} \in \{\pm 1\}^r, w > 0, \quad (14)$$

$$\mathbf{W}_2 = \mathbf{W}_1^\top. \quad (15)$$

Thus $\mathbf{W}_1 \mathbf{y}_n = dw(ex_n - p) \mathbf{w}$. Since $-p < 0$ and $e - p > 0$, corresponding to $x_n = 0$ and $x_n = 1$ respectively, we obtain $\text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) = dw((1 - x_n)p \cdot \text{ReLU}(-\mathbf{w}) + x_n(e - p) \cdot \text{ReLU}(\mathbf{w}))$. Denoting the number of ones in \mathbf{w} as β , i.e. $\beta = \sum_{i=1}^r \mathbb{1}(w_i = 1)$, we further simplify:

$$\begin{aligned} \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) &= \mathbf{W}_1^\top \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) \\ &= w^2 d ((1 - x_n)p \cdot \langle \mathbf{w}, \text{ReLU}(-\mathbf{w}) \rangle + x_n(e - p) \cdot \langle \mathbf{w}, \text{ReLU}(\mathbf{w}) \rangle) \boldsymbol{\alpha} \\ &= w^2 d ((1 - x_n)p \cdot (\beta - r) + x_n(e - p) \cdot \beta) \boldsymbol{\alpha} \\ &= w^2 d (ex_n - p) ((2\beta - r)x_n + r - \beta) \boldsymbol{\alpha}. \end{aligned}$$

Hence

$$\mathbf{z}_n = \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) = (ex_n - p) (1 + w^2 d ((2\beta - r)x_n + r - \beta)) \boldsymbol{\alpha}. \quad (16)$$

Linear. Using the fact that $\mathbf{e} = \mathbf{a} = e \cdot \boldsymbol{\alpha}$ due to weight-tying, we obtain from [Linear](#) layer that

$$\text{logit}_n = \langle \mathbf{e}, \mathbf{z}_n \rangle + b = ed(ex_n - p) (1 + w^2 d ((2\beta - r)x_n + r - \beta)) + b. \quad (17)$$

Prediction. We finally obtain that the prediction probability

$$f_\theta(x_1^n) = \sigma(\text{logit}_n) = x_n \cdot \sigma(ed(e - p) (1 + \beta w^2 d) + b) + (1 - x_n) \cdot \sigma(-edp (1 + (r - \beta)w^2 d) + b).$$

Thus we see that the prediction probability and hence the loss function $L(\cdot)$ in Eq. (11) is influenced only by the scalars e, p, w, b and β .

D. Canonical reparameterization

Building on the low-rank structure of the transformer parameters described above, we consider a special parameterization for them. A key property of this parameterization is that it covers both the global and local minima from Thm. 4 and Thm. 5 for all $(p, q) \in (0, 1)^2$. Recall that Thm. 5 characterizes local minima only for $p + q > 1$ whereas our special parameterization allows to discover local minima even for $p + q < 1$. Our construction follows the same outline as in Eqs. (12)-(17). First we start with the embedding layer.

Embedding. We let $\mathbf{e} = e \cdot \boldsymbol{\alpha}$ and $\mathbf{p}_n = -p \cdot \boldsymbol{\alpha}$ for all n where $e > 0, p = \frac{e}{2}$ and $\boldsymbol{\alpha} \in \{\pm 1\}^d / \sqrt{d}$. Thus the embedding \mathbf{x}_n from Eq. (12) simplifies to

$$\mathbf{x}_n = e \left(x_n - \frac{1}{2} \right) \boldsymbol{\alpha} \in \left\{ \pm \frac{e}{2} \right\} \boldsymbol{\alpha}. \quad (18)$$

Attention. Substituting this \mathbf{x}_n in Eq. (13), we have

$$\mathbf{y}_n = e \left(x_n - \frac{1}{2} \right) \boldsymbol{\alpha} + \langle \mathbf{v}, \boldsymbol{\alpha} \rangle \left(\sum_{i \in [n]} \text{att}_{n,i} \cdot e \left(x_i - \frac{1}{2} \right) \right) \boldsymbol{\alpha}, \quad (19)$$

where the attention weights $\text{att}_{n,i} = \frac{\exp(e^2 q^2 d^{5/2} (x_n - \frac{1}{2})(x_i - \frac{1}{2}))}{\sum_{j \in [n]} \exp(e^2 q^2 d^{5/2} (x_n - \frac{1}{2})(x_j - \frac{1}{2}))} \in (0, 1)$ for some $q \in \mathbb{R}$. Since $\langle \mathbf{v}, \boldsymbol{\alpha} \rangle \approx 0$, we let $\mathbf{v} = 0$ and obtain

$$\mathbf{y}_n = \mathbf{x}_n = e \left(x_n - \frac{1}{2} \right) \boldsymbol{\alpha}. \quad (20)$$

Feed-forward. For the feed-forward layer, we observe from Eq. (14) and Eq. (16) that for any $\mathbf{w} \in \{\pm 1\}^r$, only the number of 1's in \mathbf{w} , β , matters for the final vector \mathbf{z}_n which further interacts with the weight scalar w . Hence without loss of generality, we set \mathbf{w} to be the all-ones vector: $\mathbf{w} = \mathbf{1} \in \mathbb{R}^r$ and hence $\beta = r = 4d$. While we observe from Eq. (14) that $\mathbf{W}_2 = \mathbf{W}_1^\top$ for $p + q < 1$, we observe from the proof of the Thm. 4 for $p + q > 1$ in (Makkuva et al., 2024)-Appendix

B.2 that we need $\mathbf{W}_2 = -\mathbf{W}_1^\top$ in this scenario. Hence we consider the following parameterization that covers both these scenarios:

$$\mathbf{W}_1 = \frac{|w|}{\sqrt{d}} \mathbf{1} \cdot \boldsymbol{\alpha}^\top \in \mathbb{R}^{4d \times d}, \quad \mathbf{W}_2 = \frac{w}{\sqrt{d}} \boldsymbol{\alpha} \cdot \mathbf{1}^\top \in \mathbb{R}^{d \times 4d}.$$

Here $w > 0$ ensures $\mathbf{W}_2 = \mathbf{W}_1^\top$ whereas $w < 0$, $\mathbf{W}_2 = -\mathbf{W}_1^\top$. Using this parameterization, substituting $\beta = r = 4d$ and $w \mapsto \frac{w}{d}$ in Eq. (16), we get

$$\mathbf{z}_n = e \left(x_n - \frac{1}{2} \right) (1 + 4w|w|x_n) \boldsymbol{\alpha}. \quad (21)$$

Linear. Since $\mathbf{e} = \mathbf{a} = e \cdot \boldsymbol{\alpha}$ due to weight-tying, Eq. (17) simplifies to

$$\text{logit}_n = \langle \mathbf{e}, \mathbf{z}_n \rangle + b = e^2 \left(x_n - \frac{1}{2} \right) (1 + 4w|w|x_n) + b \quad (22)$$

$$\stackrel{(x_n = x_n^2)}{=} e^2 \left(x_n + 4w|w|x_n - \frac{1}{2} - 2w|w|x_n \right) + b \quad (23)$$

$$= e^2 (1 + 2w|w|) x_n + b - \frac{e^2}{2}. \quad (24)$$

Prediction. The next-token prediction probability is

$$f_{(\boldsymbol{\theta}, b)}(x_1^n) = \sigma \left(e^2 (1 + 2w|w|) x_n + b - \frac{e^2}{2} \right), \quad \boldsymbol{\theta} \triangleq (e, w) \in \mathbb{R}^2. \quad (25)$$

Loss. While we assumed $e > 0$ in the beginning, in view of Eq. (25) and the fact that $\boldsymbol{\alpha} \in \{\pm 1\}^d / \sqrt{d}$, we see that $e \in \mathbb{R}$ gives us the same expression for probability. Thus the final probability depends on just the three scalars $(e, w, b) \in \mathbb{R}^3$. Defining $\boldsymbol{\theta} = (e, w) \in \mathbb{R}^2$, we recall the cross-entropy loss $L(\cdot)$ from Eq. (8) in Sec. 2 for this canonical model:

$$L(\boldsymbol{\theta}, b) = -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_{n+1}} [x_{n+1} \cdot \log f_{(\boldsymbol{\theta}, b)}(x_1^n) + (1 - x_{n+1}) \cdot \log(1 - f_{(\boldsymbol{\theta}, b)}(x_1^n))]. \quad (26)$$

It turns out that we can further remove the bias b by minimizing the loss over it which we discuss in App. F. For now in the next section, we analyze when it's present as in Eq. (26).

E. Analysis of the loss with the bias, $L(\theta, b)$, in Eq. (8) and Eq. (26)

In this section, we analyze the loss function with the bias, $L(\theta, b)$, from Eq. (8) and Eq. (26), which will later be useful for studying $L(\theta)$. First we characterize the set of its critical points in \mathbb{R}^3 . To this end, we define the following sets of points

$$\Gamma_*(p, q) \triangleq \left\{ (e, w, b) \in \mathbb{R}^3 : e^2(1 + 2w|w|) = \log \frac{(1-p)(1-q)}{pq}, b - \frac{e^2}{2} = \log \frac{p}{1-p} \right\}, \quad (27)$$

$$\Gamma_{\min}(p, q) \triangleq \left\{ (e, w, b) \in \mathbb{R}^3 : e = 0, (p+q-1)(1+2w|w|) > 0, b = \log \frac{p}{q} \right\}, \quad (28)$$

$$\Gamma_{\text{sad}}(p, q) \triangleq \left\{ (e, w, b) \in \mathbb{R}^3 : e = 0, (p+q-1)(1+2w|w|) \leq 0, b = \log \frac{p}{q} \right\}. \quad (29)$$

The following result establishes that these sets exhaust all the critical points.

Theorem 6 (All critical points). *Let the input sequence be $\{x_n\}_{n=1}^N \sim (\pi, \mathbf{P})$, the transformer parameters $(\theta, b) = (e, w, b) \in \mathbb{R}^3$, and the next-token prediction loss $L(\cdot)$ be as in Eq. (26). Then all the stationary points of L are either in Γ_* , Γ_{\min} , or Γ_{sad} , i.e.*

$$\{(\theta, b) \in \mathbb{R}^3 : \nabla L(\theta, b) = 0\} = \Gamma_* \cup \Gamma_{\min} \cup \Gamma_{\text{sad}}. \quad (30)$$

Proof. We refer to App. I.1. □

Recall the definitions of local minima & maxima, global minima, and that of all the saddle points from Sec. 3.1. We are now ready to present the main result about the loss landscape of $L(\cdot)$.

Theorem 7 (Loss landscape with bias). *Let the input sequence be $\{x_n\}_{n=1}^N \sim (\pi, \mathbf{P})$, the transformer parameters $(e, w, b) \in \mathbb{R}^3$, and the next-token prediction loss $L(\cdot)$ be as in Eq. (26). Then for any $(p, q) \in (0, 1)^2$ with $p+q \neq 1$ and $N \in \mathbb{N}$,*

- (i) *the set of all global minima of L is given by $\Gamma_*(p, q)$,*
- (ii) *the set of all bad local minima of L is given by $\Gamma_{\min}(p, q)$,*
- (iii) *and the set of all saddle points of L is $\Gamma_{\text{sad}}(p, q)$.*

Furthermore, for any $\gamma_* \in \Gamma_*$, $\gamma_{\min} \in \Gamma_{\min}$, and $\gamma_{\text{sad}} \in \Gamma_{\text{sad}}$, the losses are ordered as

$$H(x_{n+1} | x_n) = L(\gamma_*) < L(\gamma_{\min}) = L(\gamma_{\text{sad}}) = H(x_{n+1}).$$

Remark 1. Note that a bad local minimum is a local minimum whose loss value is strictly less than that of the global minimum, as is the case here. Interestingly, Thm. 7 highlights that all local minima for the loss L are indeed bad local minima.

Proof. We refer to App. I.2. □

E.1. Technical lemmas

The proofs of both Thm. 6 and Thm. 7 rely on few key lemmas that we present below. First we start with the result that rewrites the loss $L(\theta, b)$ from Eq. (26) in a compact manner using the logistic function $\ell_{\log}(\cdot)$.

Lemma 2 (Loss as a logistic function). *The next-token prediction loss $L(\cdot)$ in Eq. (26) can be written as*

$$\begin{aligned} L(\theta, b) &= \frac{1}{N} \sum_{n \in [N]} \mathbb{E}[\ell_{\log}((2x_{n+1} - 1) \cdot \text{logit}_n)] \\ &= \mathbb{E}_{X, Y} \left[\ell_{\log} \left((2Y - 1) \left(e^2(1 + 2w|w|)X + b - \frac{e^2}{2} \right) \right) \right], \end{aligned} \quad (31)$$

where $(X, Y) \in \{0, 1\}^2$ are distributed according to $(X, Y) \sim (\pi, \mathbf{P})$, i.e. X is a Bernoulli random variable with $X \sim \pi \equiv \text{Bern}(p/(p+q))$ and $Y|X \sim \mathbf{P}(p, q)$, the Markov kernel.

The following lemma establishes the gradients of the loss function with respect to the parameters e , w , and b .

Lemma 3 (Gradient computation). *For any $(e, w, b) \in \mathbb{R}^3$ and the next-token prediction loss $L(\cdot)$ in Eq. (26), the gradients are given by*

$$\begin{aligned}\frac{\partial L}{\partial e} &= \mathbb{E}_X [(f_1 X + f_2)(2X(1 + 2w|w|) - 1)] \cdot e, \\ \frac{\partial L}{\partial w} &= \mathbb{E}_X [(f_1 X + f_2)X] \cdot 4e^2|w|, \\ \frac{\partial L}{\partial b} &= \mathbb{E}_X [f_1 X + f_2],\end{aligned}$$

where $X \in \{0, 1\}$ is a Bernoulli random variable with $X \sim \text{Bern}(p/(p+q))$, $f_1 = \sigma\left(2e^2w|w| + b + \frac{e^2}{2}\right) + q - 1 - \sigma\left(b - \frac{e^2}{2}\right) + p$, and $f_2 = \sigma\left(b - \frac{e^2}{2}\right) - p$.

Remark 2. It is interesting to note that the gradients for both e and w are product of an expectation term and an e factor. Also, except for scaling factors in terms of (e, w, b) , all the gradients are governed by the two expectation terms $\mathbb{E}[(f_1 X + f_2)X]$ and $\mathbb{E}[f_1 X + f_2]$. This observation plays a key role in obtaining an ordinary differential equation which yields the energy function \mathcal{E} , defined in Eq. (4).

Now we characterize the Hessian at both local-minima and saddle points.

Lemma 4 (Hessian at local-minima and saddle points). *For the canonical parameterization $\gamma = (b, e, w) \in \mathbb{R}^3$ and the next-token prediction loss $L(\cdot)$ in Eq. (26), the Hessian at any $\gamma_{\min} \in \Gamma_{\min}$ or $\gamma_{\text{sad}} \in \Gamma_{\text{sad}}$ is given by*

$$\nabla^2 L(\gamma) \Big|_{\gamma=\gamma_{\min}, \gamma_{\text{sad}}} = \pi_0 \pi_1 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2(p+q-1)(1+2w|w|) & 0 \\ 0 & 0 & 0 \end{bmatrix},$$

where $\pi_0 = \frac{q}{p+q}$ and $\pi_1 = \frac{p}{p+q}$.

Remark 3. We note that the Hessian is computed with the parameter ordering (b, e, w) .

The proofs of the lemmas are presented in App. J.

1045 **F. Analysis of the loss without bias, $L(\theta)$, and proof of Thm. 1**

 1046 The proof of Thm. 1, concerning the loss $L(\theta)$ in Eq. (9), is similar to that of Thm. 7 which studies the loss $L(\theta, b)$ with the
 1047 bias present. The main idea is to establish the analogous set of lemmas, as in App. E, when the bias is substituted with its
 1048 optimal choice. First we recall the loss function
 1049

1050
$$L(\theta) \triangleq L(\theta, b_*) = -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [x_{n+1} \cdot \log f_{(\theta, b_*)}(x_1^n) + (1 - x_{n+1}) \cdot \log(1 - f_{(\theta, b_*)}(x_1^n))], \quad (32)$$
 1051
 1052
$$b_* = \operatorname{argmin}_{b \in \mathbb{R}} L(\theta, b).$$
 1053
 1054

 1055 We start with the result that establishes a closed form expression for b_* .
 1056

 1057 **Lemma 5** (Optimal bias). *For $\theta = (e, w) \in \mathbb{R}^2$ and $b \in \mathbb{R}$, let $L(\theta, b)$ be the next-token prediction loss defined in Eq. (26).
 1058 Then, for any $\theta \in \mathbb{R}^2$, $L(\theta, b)$ is convex in b and the minimizer $b_* \triangleq \operatorname{argmin}_{b \in \mathbb{R}} L(\theta, b)$ is given by*
 1059

1060
$$\exp\left(b_* - \frac{e^2}{2}\right) = \frac{1}{2A} \left[\frac{p}{q} - 1 + \sqrt{\left(\frac{p}{q} - 1\right)^2 + 4 \cdot \frac{p}{q} \cdot A} \right], \quad A \triangleq \exp(e^2(1 + 2w|w|)). \quad (33)$$
 1061
 1062

 1063 Consequently, if $e^2(1 + 2w|w|) = \log \frac{(1-p)(1-q)}{pq}$, then $b_* - \frac{e^2}{2} = \log \frac{p}{1-p}$. If $e = 0$, then $b_* = \log \frac{p}{q}$.
 1064
 1065

 1066 Now we rewrite the loss in terms of the logistic function.
 1067

 1068 **Lemma 6** (Loss as a logistic function). *For any $\theta \in \mathbb{R}^2$, the next-token prediction loss $L(\theta)$ in Eq. (32) can be written as*
 1069

1070
$$L(\theta) = \frac{1}{N} \sum_{n \in [N]} \mathbb{E}[\ell_{\log}((2x_{n+1} - 1) \cdot \operatorname{logit}_n)]$$
 1071
 1072
$$= \mathbb{E}_{X, Y} \left[\ell_{\log} \left((2Y - 1) \left(e^2(1 + 2w|w|)X + b_* - \frac{e^2}{2} \right) \right) \right]. \quad (34)$$
 1073
 1074

 1075 where b_* follows from Eq. (33), $(X, Y) \in \{0, 1\}^2$ are distributed according to $(X, Y) \sim (\pi, \mathbf{P})$, i.e. X is a Bernoulli
 1076 random variable with $X \sim \pi \equiv \operatorname{Bern}(p/(p+q))$ and $Y|X \sim \mathbf{P}(p, q)$, the Markov kernel.
 1077

 1078 The following lemma establishes the gradients of the loss.
 1079

 1080 **Lemma 7** (Gradient computation). *For any $\theta = (e, w) \in \mathbb{R}^2$ and the next-token prediction loss $L(\theta)$ in Eq. (32), the
 1081 gradients are given by*
 1082

1083
$$\frac{\partial L}{\partial e} = \mathbb{E}_X [(f_1 X + f_2)X] \cdot 2(1 + 2w|w|)e,$$
 1084
 1085
$$\frac{\partial L}{\partial w} = \mathbb{E}_X [(f_1 X + f_2)X] \cdot 4e^2|w|,$$
 1086

 1087 where $X \in \{0, 1\}$ is a Bernoulli random variable with $X \sim \operatorname{Bern}(p/(p+q))$, $f_1 = \sigma\left(2e^2w|w| + b_* + \frac{e^2}{2}\right) + q - 1 -$
 1088 $\sigma\left(b_* - \frac{e^2}{2}\right) + p$, and $f_2 = \sigma\left(b_* - \frac{e^2}{2}\right) - p$. Further, $\pi_1 f_1 + f_2 = 0$.
 1089

 1090 **Remark 4.** We observe above that the gradients for both e and w are proportional to each other, except for the scaling
 1091 factors in terms of e and w . This forms the basis for the derivation of the energy function discussed in App. G.
 1092

 1093 The following lemma characterizes the Hessian.
 1094

 1095 **Lemma 8** (Hessian computation). *Let $\gamma = (b, \theta) \in \mathbb{R}^3$ with $\theta = (e, w) \in \mathbb{R}^2$, and $L(\gamma)$ be the next-token prediction loss
 1096 in Eq. (26) and $L(\theta)$ be the one in Eq. (32). Let the Hessian of L at γ be*
 1097

1098
$$H(\gamma) \triangleq \nabla_{\gamma}^2 L = \begin{bmatrix} H_{bb} & H_{b\theta} \\ H_{b\theta}^\top & H_{\theta\theta} \end{bmatrix} = \begin{bmatrix} \nabla_{bb}^2 L & \nabla_{b\theta}^2 L \\ (\nabla_{b\theta}^2 L)^\top & \nabla_{\theta\theta}^2 L \end{bmatrix} \in \mathbb{R}^{3 \times 3}.$$
 1099

Then the Hessian of L at $\theta \in \mathbb{R}^2$ is given by

$$H(\theta) \triangleq \nabla_{\theta\theta}^2 L = H_{\theta\theta} - H_{b\theta}^\top \cdot H_{bb}^{-1} \cdot H_{b\theta}. \quad (35)$$

Consequently, for any $\gamma = (b, e, w) \in \Gamma_{\min} \cup \Gamma_{\text{sad}}$, the Hessian $H(\theta)$ at $\theta = (e, w)$ is given by

$$H(\theta) = \pi_0 \pi_1 \begin{bmatrix} 2(p+q-1)(1+2w|w|) & 0 \\ 0 & 0 \end{bmatrix}, \quad (36)$$

where $\pi_0 = \frac{q}{p+q}$ and $\pi_1 = \frac{p}{p+q}$.

The proofs of the above lemmas are deferred to App. K. We are now ready to present the proof of Thm. 1.

F.1. Proof of Thm. 1

Proof. Let $\theta \in \mathbb{R}^2$ and $\gamma(\theta) = (\theta, b_*(\theta)) \in \mathbb{R}^3$ be its embedding in \mathbb{R}^3 with the optimal bias $b_*(\theta) = \operatorname{argmin}_{b \in \mathbb{R}} L(\theta, b)$ from Lemma 5. Define the following four sets of points:

$$\begin{aligned} \Theta_*(p, q) &\triangleq \left\{ (e, w) \in \mathbb{R}^2 : e^2(1+2w|w|) = \log \frac{(1-p)(1-q)}{pq} \right\}, \\ \Theta_{\min}(p, q) &\triangleq \left\{ (e, w) \in \mathbb{R}^2 : e = 0, (p+q-1)(1+2w|w|) > 0 \right\}, \\ \Theta_{\max}(p, q) &\triangleq \left\{ (e, w) \in \mathbb{R}^2 : e = 0, (p+q-1)(1+2w|w|) < 0 \right\}, \\ \Theta_{\text{sad}}(p, q) &\triangleq \left\{ (e, w) : e = 0, w = -1/\sqrt{2} \right\}. \end{aligned}$$

First we show that any critical point of $L : \mathbb{R}^2 \rightarrow \mathbb{R}$ has to lie in one of these sets. Then we characterize that they correspond to the set of all global minima, local minima & maxima, and saddle points respectively.

(i) Set of all critical points: Recall from Thm. 7 that for any critical point $\gamma = (\theta, b) = (e, w, b) \in \mathbb{R}^3$ of L , $\gamma \in \Gamma_* \cup \Gamma_{\min} \cup \Gamma_{\text{sad}}$. Here the main observation is that all these critical points are of the form $(\theta, b_*(\theta))$ where $\theta \in \Theta_* \cup \Theta_{\min} \cup \Theta_{\max} \cup \Theta_{\text{sad}}$. To see this, let $\gamma \in \Gamma_*$. Here we have $e^2(1+2w|w|) = \log \frac{(1-p)(1-q)}{pq}$ from Eq. (27) and hence $\theta \in \Theta_*$. Further, by Lemma 5, we have that the optimal bias for this θ satisfies $b_* - \frac{e^2}{2} = \log \frac{p}{1-p}$, which is precisely the characterization of the bias b for $\gamma = (e, w, b)$ in Eq. (27). Likewise, if $\gamma \in \Gamma_{\min} \cup \Gamma_{\text{sad}}$, we have $e = 0$ and hence $\theta \in \Theta_{\min} \cup \Theta_{\max} \cup \Theta_{\text{sad}}$. Hence by Lemma 5, $b_* = \log \frac{p}{q}$, matching that of Eq. (28) and Eq. (29). Thus the set of all critical points of L in \mathbb{R}^3 are of the form $(\theta, b_*(\theta))$ with where $\theta \in \Theta_* \cup \Theta_{\min} \cup \Theta_{\max} \cup \Theta_{\text{sad}}$. Since $\Gamma_* \cup \Gamma_{\min} \cup \Gamma_{\text{sad}}$ covers the entirety of stationary points of L in \mathbb{R}^3 , it follows that the set of all stationary points in \mathbb{R}^2 is precisely $\Theta_* \cup \Theta_{\min} \cup \Theta_{\max} \cup \Theta_{\text{sad}}$. Also, the ordering of losses directly follows from the aforementioned observation.

Now we characterize these critical points in terms of the extrema.

(ii) Set of global and local minima: From Eq. (27), for any global minimum $\gamma_* = (\theta_*, b_*(\theta_*))$ of L in \mathbb{R}^3 , we have $\theta_* \in \Theta_* \subseteq \mathbb{R}^2$. Hence by definition, Θ_* is the set of all global minima in \mathbb{R}^2 . A similar argument holds for Θ_{\min} , which establishes that it is a set of all local minima.

(iii) Set of local maxima and saddle points: From Eq. (29), for any saddle point $\gamma = (e, w, b_*(e, w))$ of L in \mathbb{R}^3 , we have that $e = 0$ and $(p+q-1)(1+2w|w|) \leq 0$. Hence $\theta = (e, w) \in \Theta_{\max} \cup \Theta_{\text{sad}}$. Suppose $\theta \in \Theta_{\max}$ which implies $e = 0, (p+q-1)(1+2w|w|) < 0$. By Lemma 8, the Hessian at θ (upto a positive scale) is a diagonal matrix with the entries $(p+q-1)(1+2w|w|) < 0$ and 0, corresponding to the directions of e and w respectively. Though one of the eigenvalue here is zero, using a continuity argument as in the proof of Thm. 7 for local minima, we can establish that θ is indeed a local maximum. Thus Θ_{\max} is a set of local minima.

Now suppose $(e, w) \in \Theta_{\text{sad}}$. Thus $e = 0$ and $w = -\frac{1}{\sqrt{2}}$. Since it lies at the intersection of Θ_{\min} and Θ_{\max} , using a neighborhood argument, it's straightforward to see that Θ_{sad} is indeed a set of saddle points.

Finally it follows that $\Theta_{\min}, \Theta_{\max}, \Theta_{\text{sad}}$ are the only set of local minima, maxima, and saddle points from the above fact about the characterization of the set of all critical points in terms of these sets and Θ_* , the ordering of the losses, and using the same argument as in the final steps of the proof of Thm. 7 with the bias. This concludes the proof. \square

G. Proofs of Thm. 2 and Thm. 8

Before we present the proofs of Thm. 2 and Thm. 8, we present few important lemmas that characterize the dynamics of the gradient flow trajectories. Recall from Sec. 4 that the trajectory $(\theta_t)_{t \geq 0}$ is governed by

$$\frac{d\theta_t}{dt} = -\nabla L(\theta_t), \quad \theta_t = (e_t, w_t) \in \mathbb{R}^2, t \geq 0, \quad (\text{GF})$$

starting with a randomly initialized θ_0 . The energy function $\mathcal{E}(\cdot, \cdot)$ is defined as

$$\mathcal{E}(e, w) \triangleq e^2 - (w^2 + \text{sign}(w) \cdot \log |w|), \quad \forall (e, w) \in \mathbb{R}^2 \setminus \text{e-axis}, \quad (37)$$

where e-axis $\triangleq \{(e, w = 0)\}$ and w-axis $\triangleq \{(e = 0, w)\}$. Note that $\mathcal{E}_{\text{sad}} = \mathcal{E}(0, -\frac{1}{\sqrt{2}}) = -\frac{1+\log 2}{2}$. We re-present the Lemma 1 from Sec. 4 below for the sake of completeness.

Lemma 9 (Constant energy along the flow). *For any $(p, q) \in (0, 1)^2$ and initialization $\theta_0 = (e_0, w_0) \in \mathbb{R}^2$, let $(\theta_t)_{t \geq 0}$ be the corresponding GF trajectory starting from θ_0 . If $w_0 \neq 0$, then the energy stays constant along the trajectory, i.e.*

$$\mathcal{E}(\theta_t) = e_t^2 - (w_t^2 + \text{sign}(w_t) \cdot \log |w_t|) = \mathcal{E}(\theta_0), \quad \forall t \geq 0. \quad (38)$$

On the other hand, if $w_0 = 0$, $w_t = 0$ for all $t \geq 0$. Hence, if we initialize on e-axis the trajectory always stays on the e-axis.

Now we establish that the GF trajectories always converge.

Lemma 10 (GF convergence). *Let $(\theta_t)_{t \geq 0}$ be a continuously differentiable GF trajectory starting from θ_0 . Then for all initializations $\theta_0 \in \mathbb{R}^2$,*

- (i) $(\theta_t)_{t \geq 0}$ is bounded,
- (ii) there exists a $\theta_{\text{lim}} \in \mathbb{R}^2$ such that $\lim_{t \rightarrow \infty} \theta_t = \theta_{\text{lim}}$ and
- (iii) $\lim_{t \rightarrow \infty} \|\nabla L(\theta_t)\| = \|\nabla L(\theta_{\text{lim}})\| = 0$.

Hence θ_{lim} is a critical point of L .

The following result characterizes the energy of the limit point.

Lemma 11 (Energy at the limit point). *Consider the same setting as in Lemma 10. If $\theta_0 \in \mathbb{R}^2 \setminus \text{e-axis}$, then $\mathcal{E}(\theta_{\text{lim}}) = \mathcal{E}(\theta_0)$. Hence θ_{lim} lies at the intersection of the contour line $\mathcal{E}(e, w) = \mathcal{E}_0$ with the set of critical points of L in \mathbb{R}^2 .*

On the other hand, if $\theta_0 \in \text{e-axis}$, then $\theta_{\text{lim}} \in \text{e-axis}$.

We now study the energy function on the w-axis which plays a key role in the GF analysis.

Lemma 12 (Analysis of the energy function). *Let $\mathcal{E}(\cdot, \cdot)$ be the energy function defined in Eq. (51) and $f(w) \triangleq \mathcal{E}(e = 0, w) = -(w^2 + \text{sign}(w) \cdot \log |w|)$ be the energy evaluated on w-axis for $w \in \mathbb{R} \setminus \{0\}$. Then*

- (i) $f : (-\infty, -1/\sqrt{2}] \rightarrow (-\infty, \mathcal{E}_{\text{sad}}]$ is monotonically increasing with $\lim_{w \rightarrow -\infty} f(w) = -\infty$ and the maximum being $f(-1/\sqrt{2}) = \mathcal{E}_{\text{sad}}$,
- (ii) $f : [-1/\sqrt{2}, 0) \rightarrow [\mathcal{E}_{\text{sad}}, -\infty)$ is monotonically decreasing with $\lim_{w \rightarrow 0^-} f(w) = -\infty$,
- (iii) $f'(-\frac{1}{\sqrt{2}}) = 0$, and
- (iv) $f : (0, \infty) \rightarrow (-\infty, \infty)$ is monotonically decreasing with $\lim_{w \rightarrow 0^+} f(w) = \infty$ and $\lim_{w \rightarrow \infty} f(w) = -\infty$.

We are now ready to prove Thm. 2 corresponding to $p + q > 1$.

G.1. Proof of Thm. 2

Proof. Let $\theta_0 = (e_0, w_0) \in \mathbb{R}^2$ be the initialization for the GF trajectory $(\theta_t)_{t \geq 0}$. Recall that

$$\begin{aligned} \mathcal{I}_{\text{min}} \triangleq & \left\{ (e, w) : w \in (-1/\sqrt{2}, 0), e \in (-g(w), g(w)), g(w) = \sqrt{w^2 - \log(-w) + \mathcal{E}_{\text{sad}}} \right\} \\ & \cup \{(e, w) : w > 0\} \cup \{(e, w) : w = 0\}, \end{aligned}$$

$$\begin{aligned} \mathcal{I}_{\text{sad}} &\triangleq \left\{ (e, w) : w \in [-1/\sqrt{2}, 0), e = \pm\sqrt{w^2 - \log(-w) + \mathcal{E}_{\text{sad}}} \right\}, \\ \mathcal{I}_{\text{max}} &\triangleq \left\{ (e, w) : e = 0, w < -1/\sqrt{2} \right\}, \\ \mathcal{I}_* &\triangleq \mathbb{R}^2 \setminus (\mathcal{I}_{\text{min}} \cup \mathcal{I}_{\text{sad}} \cup \mathcal{I}_{\text{max}}). \end{aligned}$$

We consider the cases $\theta_0 \in e$ -axis and $\theta_0 \in \mathbb{R}^2 \setminus e$ -axis separately. First recall from Thm. 1 and Eq. (??) that for $p + q > 1$, the loci of the global minima, $e^2(1 + 2w|w|) = \log \frac{(1-p)(1-q)}{pq} < 0$, lies entirely in the negative half-plane corresponding to $w < -\frac{1}{\sqrt{2}}$. On the other hand, all the local minima, maxima and the saddle points span the w -axis corresponding to $e = 0$.

(i) $\theta_0 \in \mathbb{R}^2 \setminus e$ -axis: Let $\mathcal{E}_0 = \mathcal{E}(\theta_0) \in \mathbb{R}$. By Lemmas. (9), (10), and (11), we have that the trajectory $(\theta_t)_{t \geq 0}$ always stays on the contour line $\mathcal{E}(e, w) = \mathcal{E}_0$ and converges to the limit θ_{lim} which is an intersection of this contour line with the set of critical points of L . Hence the crux of the proof is to establish where these intersections occur based on the initialization θ_0 and the initial energy \mathcal{E}_0 . This gives rise to the set of initializations $\mathcal{I}_{\text{min}}, \mathcal{I}_{\text{max}}, \mathcal{I}_{\text{sad}}$, and \mathcal{I}_* that correspond to the limit being a local minimum/maximum, a saddle point, or a global minimum.

We characterize them individually below starting with \mathcal{I}_{min} .

Initializations for local minima, \mathcal{I}_{min} . For $\theta_0 = (e_0, w_0) \in \mathbb{R}^2 \setminus e$ -axis, assume that $w_0 > 0$. Since $\mathcal{E}_0 \in \mathbb{R}$, there exists an unique $w_* > 0$ such that $f(w_*) = \mathcal{E}(0, w_*) = \mathcal{E}_0$ by Lemma 12, (iv). Further using the fact that the energy contour lines do not cross each other (by definition of a contour line) and the fact they do not intersect the e -axis (it's an energy barrier as discussed in Sec. 4), it follows that the contour line $\mathcal{E}(e, w) = \mathcal{E}_0$ stays entirely in the positive half-plane corresponding to $w > 0$ and $w_* > 0$ is the unique (and only) intersection of this line with the w -axis, and hence the set of critical points. Since the w -axis corresponding to $w > 0$ is a set of a local minima (Eq. (??)), it follows that any initialization (e_0, w_0) with $w_0 > 0$ converges to a local minimum.

Now suppose $-\frac{1}{\sqrt{2}} < w_0 < 0$ and $e_0 \in (-g(w_0), g(w_0))$, where $g(w_0) = \sqrt{w_0^2 - \log(-w_0) + \mathcal{E}_{\text{sad}}}$. Thus $|e_0| < g(w_0)$ and hence $e_0^2 - (w_0^2 - \log(-w_0)) = \mathcal{E}(e_0, w_0) = \mathcal{E}_0 < \mathcal{E}_{\text{sad}}$. Hence by Lemma 12, (iii), there is a unique intersection of the contour line $\mathcal{E}(e, w) = \mathcal{E}_0$ with the w -axis, which lies in the region $(-\frac{1}{\sqrt{2}}, 0)$. Further note that this contour line cannot intersect with the global minima loci as it lies in the half-plane $w < -\frac{1}{\sqrt{2}}$, and hence its only intersection with the set of critical points is this segment of w -axis, which is precisely the set of local minima the GF initialized on this line would converge to.

Thus we have shown that any initialization in $\mathcal{I}_{\text{min}} \setminus \cup \{(e, w) : w = 0\}$ converges to a local minimum, the set of which exhausts all the set of local minima Θ_{min} except for the origin. Below we will establish that any initialization on e -axis = $\{(e, w) : w = 0\}$ converges to the origin, implying \mathcal{I}_{min} is the full set of initializations for which the limit is a local minimum.

Initializations for saddle points, \mathcal{I}_{sad} . It's straightforward to see that for any $\theta_0 \in \mathcal{I}_{\text{sad}}$, $e_0^2 - (w_0^2 - \log(-w_0)) = \mathcal{E}(0, -\frac{1}{\sqrt{2}}) = \mathcal{E}_{\text{sad}}$. Since $-\frac{1}{\sqrt{2}} \leq w_0 < 0$, the point $(w, e) = (-\frac{1}{\sqrt{2}}, 0)$ is the only intersection of the contour line with the set of critical points, any initialization in \mathcal{I}_{sad} converges to the saddle point. On the other hand, there also exists a contour line $e_0^2 - (w_0^2 - \log(-w_0)) = \mathcal{E}_{\text{sad}}$ for $w_0 < -\frac{1}{\sqrt{2}}$ that passes through $(-\frac{1}{\sqrt{2}}, 0) \in \mathbb{R}^2$ and further intersecting with the global minima loci Θ_* . However, if we initialize on this line the flow escapes away from the saddle point and converges instead to a global minimum. To show this, it suffices to prove that $\frac{de_t}{dt} > 0$ and $\frac{dw_t}{dt} < 0$ if $e_0 > 0$ and $w_0 < -\frac{1}{\sqrt{2}}$, such that (w_0, e_0) is close to the saddle point $(-\frac{1}{\sqrt{2}}, 0)$ (the case for $e_0 < 0$ is similar as the flow is symmetric in $e \in \mathbb{R}$). From Lemma 7 and the definition of the GF, we have that

$$\begin{aligned} \frac{de_t}{dt} &= -\frac{\partial L}{\partial e}(e_0, w_0) = 2\mathbb{E}_X [(f_1 X + f_2) X] \cdot (1 - 2w_0^2)e_0 \\ \frac{dw_t}{dt} &= -\frac{\partial L}{\partial w}(e_0, w_0) = 4\mathbb{E}_X [(f_1 X + f_2) X] \cdot (-e_0^2 w_0). \end{aligned}$$

So it suffices to show that $\mathbb{E}_X [(f_1 X + f_2) X] > 0$. To establish this, we have from Lemma 5 that

$$\mathbb{E}_X [(f_1 X + f_2) X] = \mathbb{E}[X](f_1 + f_2) = \pi_1 \left(-\frac{f_2}{\pi_1} + f_2 \right) = -\pi_0 \cdot f_2.$$

1265 From the definition of f_2 and the optimal bias b_* in Lemma 7 and Lemma 5 respectively, we obtain

$$\begin{aligned}
 1266 \quad f_2 &= \sigma\left(b_* - \frac{e_0^2}{2}\right) - p = \left(1 + \exp\left(-b_* + \frac{e_0^2}{2}\right)\right)^{-1} - p \\
 1267 \quad &= \left(1 + \frac{2A}{\frac{p}{q} - 1 + \sqrt{\left(\frac{p}{q} - 1\right)^2 + 4 \cdot \frac{p}{q} \cdot A}}\right)^{-1} - p, \quad A \triangleq \exp(e_0^2(1 - 2w_0^2)). \\
 1268 \quad & \\
 1269 \quad & \\
 1270 \quad & \\
 1271 \quad & \\
 1272 \quad & \\
 1273 \quad &
 \end{aligned}$$

1274 When $e_0 = 0$, we have $A = 1$ and hence

$$1275 \quad f_2 = \left(1 + \frac{q}{p}\right)^{-1} - p = \frac{p}{p+q} - p = -\frac{p}{p+q}(p+q-1) < 0, \quad (39)$$

1276 where we used the fact that $p+q > 1$. Hence by continuity of f_2 in e_0 , for e_0 sufficiently close to 0, $f_2 < 0$ which
 1277 proves our claim about the direction of the flow close to the saddle point. By using the continuity of the flow, it follows
 1278 that GF cannot converge to saddle point when initialized on this contour line for $w_0 < -\frac{1}{\sqrt{2}}$. Thus \mathcal{I}_{sad} is the only set of
 1279 initializations for convergence to Θ_{sad} .

1280 **Initializations for local maxima, \mathcal{I}_{max} .** If $p+q > 1$, we have from Thm. 1 that $\Theta_{\text{max}} =$
 1281 $\{(e, w) \in \mathbb{R}^2 : e = 0, (1 + 2w|w|) < 0\} = \{(e, w) \in \mathbb{R}^2 : e = 0, w < -\frac{1}{\sqrt{2}}\}$. Thus for any $\theta_0 \in \Theta_{\text{max}}$, $\frac{d\theta_t}{dt} = 0$
 1282 for all $t \geq 0$ and hence $\theta_{\text{lim}} = \theta_0$. Further if we slightly perturb away from this set, from Eq. (39) it follows that the flow
 1283 diverges and hence it's an unstable set of critical points (they are local maxima indeed). Thus the only set of initializations
 1284 leading to local maxima are $\mathcal{I}_{\text{max}} = \Theta_{\text{max}}$.

1285 **Initializations for the global minima, \mathcal{I}_* .** Since the set of all critical points of L is $\Theta_* \cup \Theta_{\text{min}} \cup \Theta_{\text{max}} \cup \Theta_{\text{sad}}$, and
 1286 the initializations in \mathcal{I}_{min} , \mathcal{I}_{sad} , and \mathcal{I}_{max} converge to Θ_{min} , Θ_{sad} , and Θ_{max} respectively, it follows that the set of
 1287 initializations for which the GF converges to global minima is $\mathcal{I}_* = \mathbb{R}^2 \setminus (\mathcal{I}_{\text{min}} \cup \mathcal{I}_{\text{sad}} \cup \mathcal{I}_{\text{max}})$.

1288 In fact, since the loci of the global minima lies in the half-plane correspondint to $w < -\frac{1}{\sqrt{2}}$ when $p+q > 1$, we can
 1289 precisely determine the location of the global minimum for which the intersection occurs for any $\theta_0 \in \mathcal{I}_*$. Specifically, we
 1290 can solve the pair of equations $\mathcal{E}(e, w) = e^2 - w^2 + \log(-w) = \mathcal{E}_0$ and $e^2(1 - 2w^2) = \log\frac{(1-p)(1-q)}{pq}$ which has a unique
 1291 solution for $w < 0$ (upto a sign flip in e).

1292 **(ii) $\theta_0 \in \text{e-axis} \Rightarrow \theta_0 \in \mathcal{I}_{\text{min}}$:** If $\theta_0 = (e_0, w_0) \in \text{e-axis}$, we have that $w_0 = 0$ and hence $w_t = 0$ for all $t \geq 0$ (Lemma 9).
 1293 Lemma 10-(i) also establishes that the iterates $(\theta_t = (e_t, 0))_{t \geq 0}$ stay bounded on the e-axis and monotonically decrease.
 1294 Since the origin is the only critical point of L on the e-axis, and $\lim_{t \rightarrow \infty} \theta_t = \theta_{\text{lim}}$ exists, it follows that $\theta_{\text{lim}} = (0, 0)$, a
 1295 local minima. Thus $\theta_0 \in \mathcal{I}_{\text{min}}$.

1296 This concludes the proof for all the initializations $\theta_0 \in \mathbb{R}^2$.

1297 **Gaussian initialization $\mathcal{N}(0, \sigma^2 I_2)$.** When θ_0 is initialized according to the standard Gaussian distribution $\mathcal{N}(0, \sigma^2 I_2)$
 1298 with $\sigma^2 \ll \frac{1}{\sqrt{2}}$, we note that θ_0 lands in the set \mathcal{I}_{min} with high probability. In fact, this probability can be made arbitrarily
 1299 close to 1 depending on σ^2 . Thus this initialization will lead to a local minimum convergence on the w-axis. \square

1300 G.2. GF dynamics for $p+q < 1$

1301 **Theorem 8** (GF dynamics for $p+q < 1$). *Under the same setting as in Thm. 2 with $p+q < 1$, and any initialization*
 1302 $\theta_0 \in \mathbb{R}^2$, *the GF trajectory always converges to a $\theta_{\text{lim}} \in \mathbb{R}^2$ which is a critical point of the loss L . More specifically, θ_{lim} is*

1303 (i) *a local minimum if*

$$1304 \quad \theta_0 \in \mathcal{I}_{\text{min}} \triangleq \left\{ (e, w) : w < -1/\sqrt{2}, e \in (-g(w), g(w)), g(w) = \sqrt{w^2 - \log(-w) + \mathcal{E}_{\text{sad}}} \right\},$$

1305 (ii) *a saddle point if $\theta_0 \in \mathcal{I}_{\text{sad}} \triangleq \{(e, w) : w \leq -1/\sqrt{2}, e = \pm\sqrt{w^2 - \log(-w) + \mathcal{E}_{\text{sad}}}\}$,*

1306 (iii) *a local maximum if $\theta_0 \in \mathcal{I}_{\text{max}} \triangleq \{(e, w) : e = 0, w > -1/\sqrt{2}\}$,*

1320 (iv) and a global minimum if $\theta_0 \in \mathbb{R}^2 \setminus (\mathcal{I}_{\min} \cup \mathcal{I}_{\text{sad}} \cup \mathcal{I}_{\max})$.

1321
1322 Consequently, if we use the standard initialization $\theta_0 \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_2)$ with $\sigma^2 \ll 1/\sqrt{2}$, θ_{lim} will be a global minimum.

1323
1324 *Proof.* The proof for the case of $p + q < 1$ essentially follows the same steps as that of $p + q > 1$. If the initialization is not
1325 on the e-axis we use the energy equation to establish the convergence to the critical point at the intersection of the energy
1326 contour line with the critical set and if it starts on the e-axis, the only change is that it now converges to the global minimum
1327 instead of the origin as in the earlier case. This is due to the fact that origin turns out to be a local maximum when $p + q < 1$
1328 and hence it's an unstable critical point (which can be established as in the proof of Thm. 2 for \mathcal{I}_{\max}). \square

1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374

H. Gradient flow analysis with attention

In this section, we analyze the learning dynamics of the transformer parameters $\theta \in \mathbb{R}^3$ with the attention scalar $a \in \mathbb{R}$, i.e. $\theta = (e, w, a) \in \mathbb{R}^3$. Similar to the analysis for $\theta = (e, w) \in \mathbb{R}^2$, we first introduce the canonical parameterization including $a \in \mathbb{R}$, then analyze the corresponding loss function $L(\cdot)$ in terms of its gradients and critical points, and capitalize on it to study the gradient flow dynamics using the energy align*. We first start with the parameterization.

H.1. Canonical parameterization with attention

Embedding. Recall from App. D that we let $e = e \cdot \alpha$ and $p_n = -p \cdot \alpha$ for all n where $e > 0, p = \frac{e}{2}$ and $\alpha \in \{\pm 1\}^d / \sqrt{d}$. This results in the embedding

$$\mathbf{x}_n = e \left(x_n - \frac{1}{2} \right) \alpha.$$

Attention. Similarly, we recall from Eq. (19) that the attention output \mathbf{y}_n is given by

$$\mathbf{y}_n = e \left(x_n - \frac{1}{2} \right) \alpha + \langle \mathbf{v}, \alpha \rangle \left(\sum_{i \in [n]} \text{att}_{n,i} \cdot e \left(x_i - \frac{1}{2} \right) \right) \alpha, \quad (40)$$

where

$$\text{att}_{n,i} \triangleq \exp \left(\langle \mathbf{q}_n, \mathbf{k}_i \rangle / \sqrt{d} \right) / \left(\sum_{j \in [n]} \exp \left(\langle \mathbf{q}_n, \mathbf{k}_j \rangle / \sqrt{d} \right) \right), \quad \mathbf{q}_n = \mathbf{W}_Q \mathbf{x}_n, \quad \mathbf{k}_i = \mathbf{W}_K \mathbf{x}_i,$$

$$\mathbf{W}_Q^\top \mathbf{W}_K = (q^2 d) \alpha \cdot \alpha^\top \in \mathbb{R}^{d \times d}, \quad \text{for some } q \in \mathbb{R}.$$

Instead of the softmax, now we assume that the attention weights are linear in the scaled dot product, i.e.

$$\begin{aligned} \text{att}_{n,i} &= \frac{\langle \mathbf{q}_n, \mathbf{k}_i \rangle}{n\sqrt{d}} = \frac{1}{\sqrt{d}} \cdot \mathbf{x}_n^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_n = \frac{q^2 d}{n\sqrt{d}} \cdot (\mathbf{x}_n^\top \alpha) (\mathbf{x}_i^\top \alpha) \\ &\stackrel{(\|\alpha\|^2=d)}{=} \frac{q^2 d^3}{n\sqrt{d}} \cdot (ex_n - p)(ex_i - p) \\ &= \frac{q^2 d^{5/2}}{n} \cdot (ex_n - p)(ex_i - p) \\ &= \frac{q^2 d^{5/2} e^2}{n} \cdot \left(x_n - \frac{1}{2} \right) \left(x_i - \frac{1}{2} \right). \end{aligned} \quad (41)$$

Note that the $1/n$ factor is to ensure normalization for the attention weights in Eq. (40). Now substituting Eq. (41) in Eq. (40), we obtain

$$\begin{aligned} \mathbf{y}_n &= e \left(x_n - \frac{1}{2} \right) \alpha + \langle \mathbf{v}, \alpha \rangle \left(\sum_{i \in [n]} \text{att}_{n,i} \cdot e \left(x_i - \frac{1}{2} \right) \right) \alpha \\ &= \left[e \left(x_n - \frac{1}{2} \right) + \langle \mathbf{v}, \alpha \rangle \left(\sum_{i \in [n]} \frac{1}{n} q^2 d^{5/2} e^2 \left(x_n - \frac{1}{2} \right) \left(x_i - \frac{1}{2} \right) \right) \cdot e \left(x_i - \frac{1}{2} \right) \right] \alpha \\ &= \left[e \left(x_n - \frac{1}{2} \right) \left(1 + \langle \mathbf{v}, \alpha \rangle q^2 d^{5/2} e^2 \left(x_i - \frac{1}{2} \right)^2 \right) \right] \alpha \\ &= \left[e \left(x_n - \frac{1}{2} \right) \left(1 + \underbrace{\langle \mathbf{v}, \alpha \rangle q^2 d^{5/2} \frac{1}{4}}_a \cdot e^2 \right) \right] \alpha \\ &= e \left(x_n - \frac{1}{2} \right) (1 + ae^2) \alpha, \end{aligned}$$

1430 where we used the fact that $(x_i - \frac{1}{2})^2 = \frac{1}{4}$ since $x_i \in \{0, 1\}$, and

$$1431 \quad a \triangleq \frac{\langle \mathbf{v}, \boldsymbol{\alpha} \rangle q^2 d^{5/2}}{4} \quad (42)$$

1432 is the attention scalar. Note that this includes the scaling $\langle \mathbf{v}, \boldsymbol{\alpha} \rangle$ from the value matrix \mathbf{W}_V and q^2 from the query-key dot
1433 product. Thus we succinctly have

$$1434 \quad \mathbf{y}_n = e \left(x_n - \frac{1}{2} \right) (1 + ae^2) \boldsymbol{\alpha}. \quad (43)$$

1435 **Feed-forward.** For the feed-forward layer, we have that $\mathbf{W}_1 = \frac{|w|}{\sqrt{d}} \mathbf{1} \cdot \boldsymbol{\alpha}^\top \in \mathbb{R}^{4d \times d}$, $\mathbf{W}_2 = \frac{w}{\sqrt{d}} \boldsymbol{\alpha} \cdot \mathbf{1}^\top \in \mathbb{R}^{d \times 4d}$. Hence
1436 Eq. (43) implies

$$1437 \quad \mathbf{W}_1 \mathbf{y}_n = \frac{|w|}{\sqrt{d}} \mathbf{1} \cdot \boldsymbol{\alpha}^\top \left[e \left(x_n - \frac{1}{2} \right) (1 + ae^2) \right] \boldsymbol{\alpha} = \frac{|w|}{\sqrt{d}} \left[e \left(x_n - \frac{1}{2} \right) (1 + ae^2) \right] \mathbf{1}.$$

1438 Thus,

$$\begin{aligned} 1439 \quad \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) &= \frac{|w|}{\sqrt{d}} \mathbf{1} \cdot \text{ReLU} \left(\left[e \left(x_n - \frac{1}{2} \right) (1 + ae^2) \right] \right) \\ 1440 &= \frac{|w|}{\sqrt{d}} \mathbf{1} \cdot e \text{ReLU} \left(\left[\left(x_n - \frac{1}{2} \right) (1 + ae^2) \right] \right) \\ 1441 &= \frac{|w|}{\sqrt{d}} \mathbf{1} \cdot e \left(\frac{x_n}{2} \text{ReLU}(1 + ae^2) + \frac{1 - x_n}{2} \text{ReLU}(-1 - ae^2) \right) \\ 1442 &= \frac{|w|}{2\sqrt{d}} \mathbf{1} \cdot e \left(x_n [\text{ReLU}(1 + ae^2) - \text{ReLU}(-1 - ae^2)] + \text{ReLU}(-1 - ae^2) \right). \end{aligned}$$

1443 Using $\text{ReLU}(x) - \text{ReLU}(-x) = x$ above,

$$1444 \quad \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) = \frac{|w|}{2\sqrt{d}} \mathbf{1} \cdot e \left(x_n (1 + ae^2) + \text{ReLU}(-1 - ae^2) \right).$$

1445 Hence,

$$\begin{aligned} 1446 \quad \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) &= \frac{w}{\sqrt{d}} \boldsymbol{\alpha} \cdot \mathbf{1}^\top \frac{|w|}{2\sqrt{d}} \mathbf{1} \cdot e \left(x_n (1 + ae^2) + \text{ReLU}(-1 - ae^2) \right) \\ 1447 &= 2w|w|e \left(x_n (1 + ae^2) + \text{ReLU}(-1 - ae^2) \right) \boldsymbol{\alpha} \\ 1448 &= 2w|w|e \left(x_n (1 + ae^2) + \frac{(-1 - ae^2)}{2} + \frac{|1 + ae^2|}{2} \right) \boldsymbol{\alpha} \\ 1449 &= 2w|w|e \left(\left(x_n - \frac{1}{2} \right) (1 + ae^2) + \frac{|1 + ae^2|}{2} \right) \boldsymbol{\alpha}. \end{aligned}$$

1450 Thus the embedding \mathbf{z}_n is given by

$$\begin{aligned} 1451 \quad \mathbf{z}_n &= \mathbf{y}_n + \mathbf{W}_2 \text{ReLU}(\mathbf{W}_1 \mathbf{y}_n) = \left[e \left(x_n - \frac{1}{2} \right) (1 + ae^2) \right] \boldsymbol{\alpha} + 2w|w|e \left(\left(x_n - \frac{1}{2} \right) (1 + ae^2) + \frac{|1 + ae^2|}{2} \right) \boldsymbol{\alpha} \\ 1452 &= e \left[\left(x_n - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] \boldsymbol{\alpha}. \end{aligned}$$

1453 **Linear.** Since $\mathbf{e} = \mathbf{a} \cdot \boldsymbol{\alpha}$ due to weight-tying, the logits are given by

$$1454 \quad \text{logit}_n(e, w, a, b) = \langle \mathbf{a}, \mathbf{z}_n \rangle + b = e^2 \left[\left(x_n - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] + b.$$

1485 **Loss.** Denote $\theta \triangleq (e, w, a) \in \mathbb{R}^3$. Similar to the case without a (Eq. (11) and Lemma 2), the cross-entropy loss in our
 1486 setting can be compactly written as

$$1487$$

$$1488 \quad L(\theta, b) = \frac{1}{N} \sum_{n \in [N]} \mathbb{E}[\ell_{\log}((2x_{n+1} - 1) \cdot \text{logit}_n(\theta, b))] = \mathbb{E}_{X,Y} [\ell_{\log}((2Y - 1) \cdot \text{logit}_X(\theta, b))], \quad (44)$$

$$1489$$

$$1490$$

1491 where $\text{logit}_X(\theta, b) \triangleq e^2 \left[\left(X - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] + b$, $(X, Y) \in \{0, 1\}^2$ are distributed ac-
 1492 cording to $(X, Y) \sim (\pi, \mathbf{P})$, i.e. X is a Bernoulli random variable with $X \sim \pi \equiv \text{Bern}(p/(p+q))$ and $Y|X \sim \mathbf{P}(p, q)$,
 1493 the Markov kernel. Further, using the convexity of b in $L(\cdot, b)$, we can consider the optimal bias $b_*(\theta) = \text{argmin}_{b \in \mathbb{R}} L(\theta, b)$
 1494 in Eq. (44) to obtain the loss $L(\theta)$:

$$1495$$

$$1496 \quad L(\theta) \triangleq L(\theta, b_*) = \mathbb{E}_{X,Y} [\ell_{\log}((2Y - 1) \cdot \text{logit}_X(\theta, b_*))]$$

$$1497$$

$$1498 \quad = \mathbb{E}_{X,Y} \left[\ell_{\log} \left((2Y - 1) \cdot \left(e^2 \left[\left(X - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] + b_* \right) \right) \right]. \quad (45)$$

$$1499$$

$$1500$$

1501 H.2. Analysis of the loss function $L(\theta)$ from Eq. (45)

1502 Now we establish the gradients of the loss function.

1503 **Lemma 13** (Gradient computation and optimal bias). *For any $\theta = (e, w, a) \in \mathbb{R}^3$ and the next-token prediction loss $L(\theta)$*
 1504 *in Eq. (45), the gradients are given by*

$$1505$$

$$1506 \quad \frac{\partial L}{\partial e} = -\mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \cdot 2e (1 + ae^2) (1 + 2w|w|)$$

$$1507$$

$$1508 \quad \quad \quad - \mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \cdot 2e^3 a (1 + 2w|w|),$$

$$1509$$

$$1510 \quad \frac{\partial L}{\partial w} = -\mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \cdot 2e^2 (1 + ae^2) (|w| + \text{sign}(w) w),$$

$$1511$$

$$1512 \quad \frac{\partial L}{\partial a} = -\mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \cdot e^4 (1 + 2w|w|),$$

$$1513$$

$$1514$$

$$1515$$

1516 where $X \in \{0, 1\}$ is a Bernoulli random variable with $X \sim \text{Bern}(p/(p+q))$, and

$$1517$$

$$1518 \quad f_1 \triangleq 1 - p - q - \phi_1 + \phi_0, \quad f_2 \triangleq p - \phi_0,$$

$$1519$$

$$1520 \quad \phi_1 \triangleq \sigma \left(e^2 \left(\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right) + b_* \right),$$

$$1521$$

$$1522 \quad \phi_0 \triangleq \sigma \left(e^2 \left(\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right) + b_* \right),$$

$$1523$$

$$1524$$

1525 where the optimal bias b_* is obtained by solving $\pi_1 f_1 + f_2 = 0$.

1526 *Proof.* We defer to App. M. □

1527 **Theorem 9** (All critical points for linear attention in \mathbb{R}^4). *Let the input sequence be $\{x_n\}_{n=1}^N \sim (\pi, \mathbf{P})$, the transformer*
 1528 *parameters $\theta = (e, w, b, a) \in \mathbb{R}^4$, and the next-token prediction loss $L(\cdot)$ be as in Eq. (44). Then for any $(p, q) \in (0, 1)^2$*
 1529 *with $p + q \neq 1$ and $N \in \mathbb{N}$,*

1530 (i) *the set of all global minima is given by*

$$1531$$

$$1532 \quad \Gamma_*(p, q) \triangleq \{(e, w, b, a) \in \mathbb{R}^4 : e^2 w|w|(1 + ae^2) + b = \frac{1}{2} \log \frac{p(1-q)}{q(1-p)}, \quad (46)$$

$$1533 \quad \quad \quad e^2 (1 + ae^2) (1 + 2w|w|) = \log \frac{(1-q)(1-p)}{pq} \} \quad (47)$$

$$1534$$

$$1535$$

$$1536$$

$$1537$$

$$1538$$

$$1539$$

1540 (ii) a set of local minima is given by

$$1541 \quad \Gamma_{\min}(p, q) \triangleq \left\{ \gamma_{\min} = (e, w, b, a) \in \mathbb{R}^4 : e = 0, (p + q - 1)(1 + 2w|w|) > 0, b = \log \frac{p}{q} \right\},$$

1546 (iii) a set of saddle points is

$$1547 \quad \Gamma_{\text{sad}}(p, q) \triangleq \left\{ \gamma_{\text{sad}} = (e, w, b, a) \in \mathbb{R}^4 : e = 0, (p + q - 1)(1 + 2w|w|) \leq 0, b = \log \frac{p}{q} \right\}.$$

1550 (iv) a set of stationary points is

$$1551 \quad \Gamma_{\text{station}}(p, q) \triangleq \left\{ \gamma_{\text{station}} = (e, w, b, a) \in \mathbb{R}^4 : e \neq 0, 1 + ae^2 = 0, 1 + 2w|w| = 0, b = \log \frac{p}{q} \right\},$$

1556 Thus the set of all critical points is

$$1557 \quad \{\theta \in \mathbb{R}^2 : \nabla L(\theta) = 0\} = \Gamma_{\star} \cup \Gamma_{\min} \cup \Gamma_{\text{sad}} \cup \Gamma_{\text{station}}. \quad (48)$$

1559 In addition, for any $\theta_{\star} \in \Gamma_{\star}$, $\theta_{\min} \in \Gamma_{\min}$ and $\theta_{\text{sad}} \in \Gamma_{\text{sad}}$, the loss values satisfy

$$1561 \quad H(x_{n+1} | x_n) = L(\theta_{\star}) < L(\theta_{\min}) = L(\theta_{\max}) = L(\theta_{\text{sad}}) = H(x_{n+1}).$$

1563 **Theorem 10** (All critical points in \mathbb{R}^3). *Let the input sequence be $\{x_n\}_{n=1}^N \sim (\pi, \mathbf{P})$, the transformer parameters $\theta = (e, w, a) \in \mathbb{R}^3$, and the next-token prediction loss $L(\cdot)$ be as in Eq. (45). Then for any $(p, q) \in (0, 1)^2$ with $p + q \neq 1$ and $N \in \mathbb{N}$,*

1567 (i) the set of all global minima is given by

$$1568 \quad \Theta_{\star}(p, q) \triangleq \left\{ (e, w, a) \in \mathbb{R}^3 : e^2 (1 + ae^2) (1 + 2w|w|) = \log \frac{(1-q)(1-p)}{pq} \right\} \quad (49)$$

1571 (ii) a set of local minima is given by

$$1572 \quad \Theta_{\min}(p, q) \triangleq \left\{ (e, w, a) \in \mathbb{R}^3 : e = 0, (p + q - 1)(1 + 2w|w|) > 0 \right\},$$

1576 (iii) a set of local maxima is given by

$$1577 \quad \Theta_{\max}(p, q) \triangleq \left\{ (e, w, a) \in \mathbb{R}^3 : e = 0, (p + q - 1)(1 + 2w|w|) < 0 \right\},$$

1581 (iv) a set of saddle points is

$$1582 \quad \Theta_{\text{sad}}(p, q) \triangleq \left\{ (e, w, a) \in \mathbb{R}^3 : \left(0, -1/\sqrt{2}, a \right) \right\}.$$

1585 Defining a set of stationary points $\Theta_{\text{station}}(p, q) \triangleq \left\{ (e, w, a) \in \mathbb{R}^3 : e \neq 0, 1 + ae^2 = 0, 1 + 2w|w| = 0 \right\}$, the set of all critical points is

$$1587 \quad \{\theta \in \mathbb{R}^2 : \nabla L(\theta) = 0\} = \Theta_{\star} \cup \Theta_{\min} \cup \Theta_{\max} \cup \Theta_{\text{sad}} \cup \Theta_{\text{station}}. \quad (50)$$

1589 In addition, for any $\theta_{\star} \in \Theta_{\star}$, $\theta_{\min} \in \Theta_{\min}$, and $\theta_{\text{sad}} \in \Theta_{\text{sad}}$, the loss values satisfy

$$1591 \quad H(x_{n+1} | x_n) = L(\theta_{\star}) < L(\theta_{\min}) = L(\theta_{\max}) = L(\theta_{\text{sad}}) = H(x_{n+1}).$$

1593 We defer the proofs of the theorems to App. M.2.

1594

1595 **H.3. Gradient flow analysis**

 1596 Analogous to the gradient flow analysis for $\theta = (e, w) \in \mathbb{R}^2$ in **G**, we now study its counterpart together with the attention
 1597 scalar, i.e. $\theta = (e, w, a) \in \mathbb{R}^3$. To this end, let $(\theta_t)_{t \geq 0}$ be a C^1 curve in \mathbb{R}^3 governed by
 1598

1599
$$\frac{d\theta_t}{dt} = -\nabla L(\theta_t), \quad \theta_t = (e_t, w_t, a_t) \in \mathbb{R}^3, t \geq 0, \quad (\text{GF-Attention})$$

 1600 starting with a randomly initialized θ_0 . We define the *energy function* $\mathcal{E}(\cdot, \cdot, \cdot)$ as
 1601

1602
$$\mathcal{E}(e, w, a) \triangleq e^2 - (w^2 + \text{sign}(w) \cdot \log |w|) - 2a^2, \quad \forall (e, w, a) \in \mathbb{R}^3 \setminus \text{ea-plane}, \quad (51)$$

 1603 where ea-plane $\triangleq \{(e, w = 0, a)\}$. The following lemma presents the crucial result that the energy is constant along the
 1604 flow in **GF-Attention**.
 1605

 1606 **Lemma 14** (Constant energy along the flow). *For any $(p, q) \in (0, 1)^2$ and initialization $\theta_0 = (e_0, w_0, a_0) \in \mathbb{R}^3$, let
 1607 $(\theta_t)_{t \geq 0}$ be the corresponding **GF-Attention** trajectory starting from θ_0 . If $w_0 \neq 0$, then the energy stays constant along the
 1608 trajectory, i.e.*
 1609

1610
$$\mathcal{E}(\theta_t) = e_t^2 - (w_t^2 + \text{sign}(w_t) \cdot \log |w_t|) - 2a_t^2 = \mathcal{E}(\theta_0), \quad \forall t \geq 0. \quad (52)$$

 1611 *On the other hand, if $w_0 = 0$, $w_t = 0$ for all $t \geq 0$. Hence, if we initialize on ea-plane the trajectory always stays on the
 1612 ea-plane.*
 1613

 1614 Now we characterize the convergence of the gradient flow.
 1615

 1616 **Lemma 15** (GF convergence). *Let $(\theta_t)_{t \geq 0}$ be a continuously differentiable **GF-Attention** trajectory starting from θ_0 . Then
 1617 for all initializations $\theta_0 \in \mathbb{R}^3$,*
 1618

- 1619 (i)
- $(\theta_t)_{t \geq 0}$
- is bounded,
-
- 1620 (ii) there exists a
- $\theta_{\text{lim}} \in \mathbb{R}^3$
- such that
- $\lim_{t \rightarrow \infty} \theta_t = \theta_{\text{lim}}$
- and
-
- 1621 (iii)
- $\lim_{t \rightarrow \infty} \|\nabla L(\theta_t)\| = \|\nabla L(\theta_{\text{lim}})\| = 0$
- .
-
- 1622

 1623 Hence θ_{lim} is a critical point of L .
 1624

 1625 The following result characterizes the energy of the limit point.
 1626

 1627 **Lemma 16** (Energy at the limit point). *Consider the same setting as in Lemma 15. If $\theta_0 \in \mathbb{R}^3 \setminus \text{ea-plane}$, then
 1628 $\mathcal{E}(\theta_{\text{lim}}) = \mathcal{E}(\theta_0)$. Hence θ_{lim} lies at the intersection of the contour line $\mathcal{E}(e, w) = \mathcal{E}_0$ with the set of critical points of L in
 1629 \mathbb{R}^3 .*
 1630

 1631 *On the other hand, if $\theta_0 \in \text{ea-plane}$, then $\theta_{\text{lim}} \in \text{ea-plane}$.*
 1632

 1633 We defer the proofs of the lemmas to App. M.
 1634

 1635 **H.4. Role of Standard Initialization**

 1636 **Theorem 11** ([Informal] Role of standard initialization for $p + q - 1 > 0$). *If we use the standard initialization $\theta_0 \sim$
 1637 $\mathcal{N}(0, \sigma^2 \mathbf{I}_2)$ with $\sigma^2 \ll 1/\sqrt{2}$, θ_{lim} will be a local minimum with high probability.*
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649

I. Proofs of theorems in App. E

I.1. Proof of Thm. 7

Proof. We characterize the set of global minima, local minima, and that of the saddle points individually.

(i) Set of all global minima. Let $\gamma_\star \in \mathbb{R}^3$ be arbitrary. From (?)Lemma 1]makkuva2024attention, we have that γ_\star is a global minimum for the loss $L(\cdot)$ in Eq. (26) if and only if its prediction probability satisfies $f_{\gamma_\star}(x_1^n) = \mathbb{P}(x_{n+1} = 1 | x_n)$, the Markov kernel. Since the input $\{x_n\}_{n=1}^N \sim (\pi(p, q), \mathbf{P}(p, q))$, we have that

$$\mathbb{P}(x_{n+1} = 1 | x_n) = (1 - x_n)p + x_n(1 - q) = (1 - p - q)x_n + p. \quad (53)$$

On the other hand, by definition, from Eq. (7), $f_{\gamma_\star}(x_1^n) = \sigma\left(e^2(1 + 2w|w|)x_n + b - \frac{e^2}{2}\right)$, where $\gamma_\star = (e, w, b)$. Since $x_n \in \{0, 1\}$, this can be further simplified to

$$\begin{aligned} f_{\gamma_\star}(x_1^n) &= \sigma\left(e^2(1 + 2w|w|)x_n + b - \frac{e^2}{2}\right) \\ &= x_n \cdot \sigma\left(e^2(1 + 2w|w|) + b - \frac{e^2}{2}\right) + (1 - x_n) \cdot \sigma\left(b - \frac{e^2}{2}\right) \\ &= x_n \left(\sigma\left(2e^2w|w| + b + \frac{e^2}{2}\right) - \sigma\left(b - \frac{e^2}{2}\right)\right) + \sigma\left(b - \frac{e^2}{2}\right). \end{aligned} \quad (54)$$

Since both $f_{\gamma_\star}(x_1^n)$ and $\mathbb{P}(x_{n+1} = 1 | x_n)$ are linear functions of x_n , equating them for all values of $x_n \in \{0, 1\}$ implies that the respective coefficients in these functions in Eq. (53) and Eq. (54) are also equal, i.e.

$$\begin{aligned} \sigma\left(b - \frac{e^2}{2}\right) &= p, \\ \sigma\left(2e^2w|w| + b + \frac{e^2}{2}\right) - \sigma\left(b - \frac{e^2}{2}\right) &= 1 - p - q, \end{aligned}$$

and hence

$$\sigma\left(b - \frac{e^2}{2}\right) = p, \quad \sigma\left(2e^2w|w| + b + \frac{e^2}{2}\right) = 1 - q. \quad (55)$$

Since $\sigma(z) = y$ for $y \in (0, 1)$ implies $z = \log \frac{y}{1-y}$, Eq. (55) can be rewritten as

$$b - \frac{e^2}{2} = \log \frac{p}{1-p}, \quad 2e^2w|w| + b + \frac{e^2}{2} = \log \frac{1-q}{q}.$$

Using $2e^2w|w| + b + \frac{e^2}{2} = e^2(1 + 2w|w|) + b - \frac{e^2}{2} = e^2(1 + 2w|w|) + \log \frac{p}{1-p}$ in the second equality above, we obtain

$$\begin{aligned} b - \frac{e^2}{2} &= \log \frac{p}{1-p}, \\ e^2(1 + 2w|w|) &= \log \frac{1-q}{q} + \log \frac{1-p}{p} = \log \frac{(1-p)(1-q)}{pq}. \end{aligned} \quad (56)$$

Thus $\gamma_\star \in \mathbb{R}^3$ is a global minimum for $L(\cdot)$ if and only if it satisfies Eq. (56) (note that it's already a critical point, as established in Thm. 6). Thus, the set of all global minimum $\Gamma_\star(p, q)$ is given by

$$\Gamma_\star(p, q) \triangleq \left\{ \gamma_\star = (e, w, b) \in \mathbb{R}^3 : e^2(1 + 2w|w|) = \log \frac{(1-p)(1-q)}{pq}, b - \frac{e^2}{2} = \log \frac{p}{1-p} \right\}.$$

Since the prediction $f_{\gamma_\star}(\cdot)$ equals the Markov kernel for any $\gamma_\star \in \Gamma_\star$, it follows from Thm. 4 (or (?)Lemma 1]makkuva2024attention) that $L(\gamma_\star) = H(x_{n+1} | x_n)$, the entropy rate of the Markov chain.

(ii) Set of local minima and saddle points.

Define $\Gamma_{\min}(p, q) \subseteq \mathbb{R}^3$ and $\Gamma_{\text{sad}} \subseteq \mathbb{R}^3$ as follows:

$$\begin{aligned}\Gamma_{\min}(p, q) &\triangleq \left\{ \gamma_{\min} = (e, w, b) \in \mathbb{R}^3 : e = 0, (p + q - 1)(1 + 2w|w|) > 0, b = \log \frac{p}{q} \right\}, \\ \Gamma_{\text{sad}}(p, q) &\triangleq \left\{ \gamma_{\text{sad}} = (e, w, b) \in \mathbb{R}^3 : e = 0, (p + q - 1)(1 + 2w|w|) \leq 0, b = \log \frac{p}{q} \right\}.\end{aligned}$$

To show that Γ_{\min} is the set of all bad local minima for $L(\cdot)$, we first show that any $\gamma_{\min} \in \Gamma_{\min}$ is a bad local minimum and then show that every bad local minimum should belong to Γ_{\min} . Similarly for Γ_{sad} . We start with the local minima.

Let $\gamma_{\min} = (e, w, b) \in \Gamma_{\min}$. Recall that γ_{\min} is a stationary point (Thm. 6), i.e.

$$\nabla L(\gamma_{\min}) = 0.$$

Rearranging the order of scalars and writing $\gamma_{\min} = (b, e, w)$, from Lemma 4, the Hessian of the loss at γ_{\min} is

$$\nabla^2 L(\gamma_{\min}) = \pi_0 \pi_1 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2(p + q - 1)(1 + 2w|w|) & 0 \\ 0 & 0 & 0 \end{bmatrix}. \quad (57)$$

By definition, $\gamma_{\min} = (b, e, w)$ satisfies $(p + q - 1)(1 + 2w|w|) > 0$. Thus its Hessian in Eq. (57) has a block diagonal structure of the form $\begin{bmatrix} \mathbf{H}_{b,e} & 0 \\ 0 & 0 \end{bmatrix}$ where $\mathbf{H}_{b,e}$ has both the eigen values positive, and hence positive-definite. In other words, γ_{\min} is a local minimum for $L(\cdot)$ in the $(b, e) \in \mathbb{R}^2$ space for any fixed w in the set. Interestingly, using the continuity argument and the fact that $L(b = \log \frac{p}{q}, e = 0, w)$ is constant in $w \in \mathbb{R}$, we can essentially follow the same steps as in proof of Theorem 2 in (?)Appendix B.3]makkuva2024attention (Thm. 5 above) to show that $\gamma_{\min} = (b, e, w)$ is a also local minimum for $L(\cdot)$ in the full parameter space \mathbb{R}^3 . This establishes that γ_{\min} is a local minimum for $L(\cdot)$.

For the saddle points, let $\gamma_{\text{sad}} = (e, w, b) \in \Gamma_{\text{sad}}$. We have that γ_{sad} is a stationary point (Thm. 6) and Lemma 4 implies its Hessian (after rearranging the order of scalars as above with $\gamma_{\text{sad}} = (b, e, w)$) is:

$$\nabla^2 L(\gamma_{\text{sad}}) = \pi_0 \pi_1 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2(p + q - 1)(1 + 2w|w|) & 0 \\ 0 & 0 & 0 \end{bmatrix}. \quad (58)$$

If $w \neq -\frac{1}{\sqrt{2}}$, $(p + q - 1)(1 + 2w|w|) < 0$ for any $\gamma_{\text{sad}} \in \Gamma_{\text{sad}}$, and hence the Hessian $\nabla^2 L(\gamma_{\text{sad}})$ in Eq. (58) as both positive, negative, and zero eigen values. Thus γ_{sad} is a saddle point for $L(\cdot)$. Using a neighborhood argument, we can similarly argue for $w = \frac{1}{\sqrt{2}}$ to establish that it's also a saddle point. Now we compute the loss value.

For any $\gamma_{\min} = (e, w, b) \in \Gamma_{\min}$ or $\gamma_{\text{sad}} = (e, w, b) \in \Gamma_{\text{sad}}$, we have that $e = 0$ and $b = \log \frac{p}{q}$. Thus for $\gamma = \gamma_{\min}$ or Γ_{sad} , the prediction probability in view of Eq. (7) is

$$\mathbb{P}_{\gamma}(x_{n+1} = 1 \mid x_1^n) = \sigma \left(e^2(1 + 2w|w|)x_n + b - \frac{e^2}{2} \right) = \sigma(b) = \frac{p}{p+q} = \mathbb{P}(x_{n+1} = 1),$$

the marginal distribution. Substituting this equality in the definition of cross-entropy loss $L(\cdot)$ in Eq. (1) and the fact that $\mathbb{P}(x_{n+1} = 1) = \frac{p}{p+q} = \pi_1$, following the same steps as in (?)Appendix B.3]makkuva2024attention, we obtain

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

$$= H(\boldsymbol{\pi}) = H(x_{n+1}).$$

Thus $L(\boldsymbol{\gamma}_{\min}) = L(\boldsymbol{\gamma}_{\text{sad}}) = H(x_{n+1})$. To see that $H(x_{n+1} | x_n) = L(\boldsymbol{\gamma}_*) < L(\boldsymbol{\gamma}_{\min}) = L(\boldsymbol{\gamma}_{\text{sad}}) = H(x_{n+1})$ for any global minimum $\boldsymbol{\gamma}_*$, observe that the gap

$$L(\boldsymbol{\gamma}_{\min}) - L_* = H(x_{n+1}) - H(x_{n+1}|x_n) = I(x_n; x_{n+1}) \geq 0,$$

where $I(x_n; x_{n+1})$ is the mutual information between x_n and x_{n+1} (Cover and Thomas, 2006). Hence the optimality gap equals zero if and only if the mutual information equals zero, which happens when x_n and x_{n+1} are independent, i.e. $\mathbb{P}(x_{n+1} = 1 | x_n)$ doesn't depend on x_n . Since $\mathbb{P}(x_{n+1} = 1 | x_n) = (1 - p - q)x_n + p$ from Eq. (53), this happens only when $p + q = 1$ which contradicts the theorem assumption that $p + q \neq 1$. Hence $H(x_{n+1} | x_n) = L(\boldsymbol{\gamma}_*) < L(\boldsymbol{\gamma}_{\min}) = L(\boldsymbol{\gamma}_{\text{sad}}) = H(x_{n+1})$.

Now we finally show that $\boldsymbol{\Gamma}_{\min}$ and $\boldsymbol{\Gamma}_{\text{sad}}$ are the only set of bad local minima and saddle points respectively. Let $\boldsymbol{\gamma}$ is a bad local minimum for $L(\cdot)$. By definition, it's also a critical point. Recall from Thm. 6 that any stationary point $\boldsymbol{\gamma} = (e, w, b)$ for the loss $L(\cdot)$ satisfies that either $\boldsymbol{\gamma} \in \boldsymbol{\Gamma}_*$, $\boldsymbol{\gamma} \in \boldsymbol{\Gamma}_{\min}$, or $\boldsymbol{\gamma} \in \boldsymbol{\Gamma}_{\text{sad}}$. Clearly $\boldsymbol{\gamma} \notin \boldsymbol{\Gamma}_*$, as $\boldsymbol{\Gamma}_*$ is the set of all global minima. Similarly, $\boldsymbol{\gamma} \notin \boldsymbol{\Gamma}_{\text{sad}}$ as every point in $\boldsymbol{\Gamma}_{\text{sad}}$ is a saddle point for the loss $L(\cdot)$ as established above. Hence $\boldsymbol{\gamma} \in \boldsymbol{\Gamma}_{\min}$. Thus every bad local minimum in \mathbb{R}^3 belongs to $\boldsymbol{\Gamma}_{\min}$. This coupled with the fact above that $\boldsymbol{\Gamma}_{\min}$ is a set of bad local minima implies $\boldsymbol{\Gamma}_{\min}$ is indeed the set of all bad local minima. The proof for $\boldsymbol{\Gamma}_{\text{sad}}$ is similar. \square

I.2. Proof of Thm. 6

Proof. Let $\boldsymbol{\gamma} = (e, w, b) \in \mathbb{R}^3$ be such that $\nabla L(\boldsymbol{\gamma}) = \left(\frac{\partial L}{\partial e}, \frac{\partial L}{\partial w}, \frac{\partial L}{\partial b} \right)^\top = 0$. By Lemma 3, we have

$$\frac{\partial L}{\partial e} = \mathbb{E}_X [(f_1 X + f_2)(2X(1 + 2w|w|) - 1)] \cdot e = 0, \quad (59)$$

$$\frac{\partial L}{\partial w} = \mathbb{E}_X [(f_1 X + f_2)X] \cdot 4e^2|w| = 0, \quad (60)$$

$$\frac{\partial L}{\partial b} = \mathbb{E}_X [f_1 X + f_2] = 0, \quad (61)$$

where $X \sim \text{Bern}(p/(p+q))$, $f_1 = \sigma\left(2e^2 w|w| + b + \frac{e^2}{2}\right) + q - 1 - \sigma\left(b - \frac{e^2}{2}\right) + p$, and $f_2 = \sigma\left(b - \frac{e^2}{2}\right) - p$. Our goal is to now show that Eqs. (59)-(61) hold only if either $(e = 0, b = \log \frac{p}{q})$ or $(f_1 = 0, f_2 = 0)$. We consider two cases corresponding to $e = 0$ and $e \neq 0$.

(i): $e = 0$. If $e = 0$, we readily see that $\frac{\partial L}{\partial e} = \frac{\partial L}{\partial e} = 0$. Further, $f_1 = p + q - 1$ and $f_2 = \sigma(b) - p$. Hence, Eq. (61) implies that

$$\mathbb{E}_X [f_1 X + f_2] = (p + q - 1)\mathbb{E}[X] + \sigma(b) - p = (p + q - 1)\frac{p}{p+q} + \sigma(b) - p = \sigma(b) - \frac{p}{p+q} = 0,$$

which implies that $b = \log \frac{p}{q}$. Since $w \in \mathbb{R}$ is arbitrary, we see in this case that

$$\begin{aligned} \boldsymbol{\gamma} \in & \left\{ (e, w, b) \in \mathbb{R}^3 : e = 0, (p + q - 1)(1 + 2w|w|) > 0, b = \log \frac{p}{q} \right\} \cup \\ & \left\{ (e, w, b) \in \mathbb{R}^3 : e = 0, (p + q - 1)(1 + 2w|w|) < 0, b = \log \frac{p}{q} \right\}, \\ & = \boldsymbol{\Gamma}_{\min} \cup \boldsymbol{\Gamma}_{\text{sad}}. \end{aligned}$$

(ii): $e \neq 0$. Suppose $e \neq 0$. Here we show that $f_1 = f_2 = 0$ and hence $\boldsymbol{\gamma} \in \boldsymbol{\Gamma}_*$. We consider two cases corresponding to $w \neq 0$ and $w = 0$. Let $w \neq 0$. Since both $e \neq 0$ and $w \neq 0$, and $X = X^2$ in Eq. (60) with $\mathbb{E}[X] = \pi_1 = \frac{p}{p+q} > 0$, we obtain that

$$\mathbb{E}_X [f_1 X^2 + f_2 X] = (f_1 + f_2)\mathbb{E}[X] = (f_1 + f_2)\pi_1 = 0,$$

and hence $f_1 + f_2 = 0$. Further Eq. (61) implies that $f_1 \pi_1 + f_2 = 0$. Together, this implies $f_1 = f_2 = 0$.

1815 Now suppose $w = 0$. From Eq. (61), we have that $\mathbb{E}[f_1 X + f_2] = f_1 \pi_1 + f_2 = 0$. Since $e \neq 0$, Eq. (59) yields

$$1816 \mathbb{E}_X [(f_1 X + f_2)(2X - 1)] = 2\mathbb{E}_X [(f_1 X + f_2)X] = 2\pi_1(f_1 + f_2) = 0.$$

1818 So $f_1 = f_2 = 0$ in this case too. Thus we have showed that whenever $e \neq 0$, we have $f_1 = f_2 = 0$. Recalling the
1819 expressions for f_1 and f_2 ,

$$1821 f_2 = \sigma\left(b - \frac{e^2}{2}\right) - p = 0 \Rightarrow b - \frac{e^2}{2} = \log \frac{p}{1-p}, \quad (62)$$

$$1822 f_1 = \sigma\left(2e^2 w|w| + b + \frac{e^2}{2}\right) + q - 1 - \sigma\left(b - \frac{e^2}{2}\right) + p = \sigma\left(2e^2 w|w| + b + \frac{e^2}{2}\right) + q - 1 = 0,$$

1823 and hence,

$$1824 2e^2 w|w| + b + \frac{e^2}{2} = \log \frac{1-q}{q}.$$

1825 Substituting $b - \frac{e^2}{2} = \log \frac{p}{1-p}$ in the above equation,

$$1826 2e^2 w|w| + e^2 + b - \frac{e^2}{2} = 2e^2 w|w| + e^2 + \log \frac{p}{1-p} = \log \frac{1-q}{q},$$

1827 and thus,

$$1828 e^2(1 + 2w|w|) = \log \frac{(1-p)(1-q)}{pq}. \quad (63)$$

1829 In view of Eq. (62) and Eq. (63), we have that $\gamma = (e, w, b) \in \Gamma_*$.

1830 Together, we have shown that whenever $\nabla L(\gamma) = 0$, we have $\gamma \in \Gamma_* \cup \Gamma_{\min} \cup \Gamma_{\text{sad}}$. Since $\Gamma_* \cup \Gamma_{\min} \cup \Gamma_{\text{sad}} \subseteq$
1831 $\{\gamma : \nabla L(\gamma) = 0\}$, we are done. \square

1832 J. Proofs of technical lemmas in App. E

1833 J.1. Proof of Lemma 2

1834 *Proof.* Recall from Eq. (11) that the cross-entropy loss $L(\cdot)$ is defined as

$$1835 L(\gamma) = -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [x_{n+1} \cdot \log f_\gamma(x_1^n) + (1 - x_{n+1}) \cdot \log(1 - f_\gamma(x_1^n))], \quad (64)$$

1836 where $f_\gamma(x_1^n) = \sigma(\text{logit}_n) = \sigma\left(e^2(1 + 2w|w|)x_n + b - \frac{e^2}{2}\right)$ from Eq. (7). For any $Y \in \{0, 1\}, Z \in \mathbb{R}$, using the fact
1837 that $1 - \sigma(Z) = \sigma(-Z)$, we have

$$1838 -[Y \log \sigma(Z) + (1 - Y) \log(1 - \sigma(Z))] = -[Y \log \sigma(Z) + (1 - Y) \log \sigma(-Z)]$$

$$1839 \stackrel{(a)}{=} Y \log(1 + \exp(-Z)) + (1 - Y) \log(1 + \exp(Z))$$

$$1840 = \log(1 + \exp(-(2Y - 1)Z))$$

$$1841 \stackrel{(b)}{=} \ell_{\log}((2Y - 1)Z), \quad (65)$$

1842 where (a) and (b) follow from the definitions of sigmoid and the logistic functions: $\sigma(z) = \frac{1}{1 + \exp(-z)}$, $\ell_{\log}(z) =$
1843 $\log(1 + \exp(-z))$ for $z \in \mathbb{R}$. Substituting $Y = x_{n+1} \in \{0, 1\}$ and $Z = \text{logit}_n \in \mathbb{R}$ in Eq. (65), Eq. (64) simplifies to

$$1844 L(\gamma) = \frac{1}{N} \sum_{n \in [N]} \mathbb{E}[\ell_{\log}((2x_{n+1} - 1) \cdot \text{logit}_n)].$$

1870 Since logit_n is only a function of x_n , the above expectation is over the distribution of the pairs (x_n, x_{n+1}) , which for all
 1871 $n \in [N]$ have the same law as a pair of random variables (X, Y) with $X \sim \pi \equiv \text{Bern}(p/(p+q))$ and $Y|X \sim \mathbf{P}(p, q)$, the
 1872 Markov kernel. Hence the above equality can be rewritten using the definition of logit_n as

$$1873 L(\gamma) = \frac{1}{N} \sum_{n \in [N]} \mathbb{E}[\ell_{\log}((2x_{n+1} - 1) \cdot \text{logit}_n)] = \mathbb{E}_{X,Y} \left[\ell_{\log} \left((2Y - 1) \left(e^2(1 + 2w|w|)X + b - \frac{e^2}{2} \right) \right) \right].$$

1877 \square

1879 J.2. Proof of Lemma 3

1880 *Proof.* With $\gamma = (e, w, b)$ and θ denoting either of the scalars e, w , or b , we have from (Lemma 2)makkuva2024attention
 1881 that the gradient of the loss $L(\cdot)$ is given by

$$1882 \nabla_{\theta} L(\gamma) = -\frac{1}{N} \sum_{n \in [N]} \mathbb{E}_{x_1^{n+1}} [(x_{n+1} - f_{\theta}(x_1^n)) \cdot \nabla_{\theta} \text{logit}_n], \quad (66)$$

1883 where $f_{\gamma}(x_1^n) = \sigma(\text{logit}_n) = \sigma \left(e^2(1 + 2w|w|)x_n + b - \frac{e^2}{2} \right)$. Using the same argument as in the proof of Lemma 6, we
 1884 can replace the expectations in Eq. (66) with that of a pair of random variables (X, Y) with $X \sim \pi \equiv \text{Bern}(p/(p+q))$
 1885 and $Y|X \sim \mathbf{P}(p, q)$, the Markov kernel. That is,

$$1886 \nabla_{\theta} L(\gamma) = -\mathbb{E}_{X,Y} \left[\left(Y - \sigma \left(e^2(1 + 2w|w|)X + b - \frac{e^2}{2} \right) \right) \cdot \nabla_{\theta} \left(e^2(1 + 2w|w|)X + b - \frac{e^2}{2} \right) \right]. \quad (67)$$

1887 Now we define the error term $\mathcal{E}(X, Y) \triangleq - \left(Y - \sigma \left(e^2(1 + 2w|w|)X + b - \frac{e^2}{2} \right) \right)$. Our goal is to show that $\mathbb{E}[\mathcal{E}(X, Y) |$
 1888 $X] = f_1 X + f_2$, where $f_1 \triangleq \sigma \left(2e^2 w|w| + b + \frac{e^2}{2} \right) + q - 1 - \sigma \left(b - \frac{e^2}{2} \right) + p$, and $f_2 \triangleq \sigma \left(b - \frac{e^2}{2} \right) - p$, which suffices
 1889 to prove the lemma. To this end, using the fact that $X \in \{0, 1\}$, we have

$$1890 \begin{aligned} \mathcal{E}(X, Y) &= - \left(Y - \sigma \left(e^2(1 + 2w|w|)X + b - \frac{e^2}{2} \right) \right) \\ 1901 &= - \left(Y - X \cdot \sigma \left(2e^2 w|w| + b + \frac{e^2}{2} \right) - (1 - X) \cdot \sigma \left(b - \frac{e^2}{2} \right) \right) \\ 1902 &= -Y + X \left(\sigma \left(2e^2 w|w| + b + \frac{e^2}{2} \right) - \sigma \left(b - \frac{e^2}{2} \right) \right) + \sigma \left(b - \frac{e^2}{2} \right). \end{aligned}$$

1903 Now taking the conditional expectation with respect to X and using the fact that $\mathbb{E}[Y|X] = \mathbb{P}(Y = 1 | X) = (1-p-q)X + p$
 1904 (since $Y|X \sim \mathbf{P}(p, q)$), we have

$$1905 \begin{aligned} \mathbb{E}[\mathcal{E}(X, Y) | X] &= -(1-p-q)X - p + X \left(\sigma \left(2e^2 w|w| + b + \frac{e^2}{2} \right) - \sigma \left(b - \frac{e^2}{2} \right) \right) + \sigma \left(b - \frac{e^2}{2} \right) \\ 1910 &= X \left(\sigma \left(2e^2 w|w| + b + \frac{e^2}{2} \right) + q - 1 - \sigma \left(b - \frac{e^2}{2} \right) + p \right) + \sigma \left(b - \frac{e^2}{2} \right) - p \\ 1911 &\stackrel{(a)}{=} f_1 X + f_2, \end{aligned}$$

1912 where (a) follows from the definition of f_1 and f_2 above. Thus Eq. (67) simplifies to

$$1913 \nabla_{\theta} L(\gamma) = \mathbb{E}_X \left[(f_1 X + f_2) \cdot \nabla_{\theta} \left(e^2(1 + 2w|w|)X + b - \frac{e^2}{2} \right) \right].$$

1914 Letting $\theta = e, w$, and b in the above equation, we finally obtain the individual gradients:

$$1915 \frac{\partial L}{\partial e} = \mathbb{E}_X [(f_1 X + f_2)(2X(1 + 2w|w| - 1))] \cdot e,$$

$$\begin{aligned}\frac{\partial L}{\partial w} &= \mathbb{E}_X [(f_1 X + f_2) X] \cdot 4e^2 |w|, \\ \frac{\partial L}{\partial b} &= \mathbb{E}_X [f_1 X + f_2].\end{aligned}$$

□

J.3. Proof of Lemma 4

Proof. Slightly changing the variable order, for any $\gamma = (b, e, w) \in \mathbb{R}^3$, we define

$$\mathbf{H}(\gamma) \triangleq \nabla^2 L(\gamma) = \begin{bmatrix} \frac{\partial^2 L}{\partial b^2} & \frac{\partial^2 L}{\partial b \partial e} & \frac{\partial^2 L}{\partial b \partial w} \\ \frac{\partial^2 L}{\partial e \partial b} & \frac{\partial^2 L}{\partial e^2} & \frac{\partial^2 L}{\partial e \partial w} \\ \frac{\partial^2 L}{\partial w \partial b} & \frac{\partial^2 L}{\partial w \partial e} & \frac{\partial^2 L}{\partial w^2} \end{bmatrix} \in \mathbb{R}^{3 \times 3}. \quad (68)$$

Recall that for any $\gamma_{\min} \in \mathbf{\Gamma}_{\min}$ and $\gamma_{\text{sad}} \in \mathbf{\Gamma}_{\text{sad}}$, we have $e = 0$ and $b = \log \frac{p}{q}$. Now we compute the second derivatives of L with respect to any $\gamma = (b = \log \frac{p}{q}, e = 0, w)$. We start with the first derivatives. By Lemma 7, the gradients are

$$\begin{aligned}\frac{\partial L}{\partial b} &= \mathbb{E}_X [f_1 X + f_2], \\ \frac{\partial L}{\partial e} &= \mathbb{E}_X [(f_1 X + f_2)(2X(1 + 2w|w|) - 1)] \cdot e, \\ \frac{\partial L}{\partial w} &= \mathbb{E}_X [(f_1 X + f_2) X] \cdot 4e^2 |w|,\end{aligned} \quad (69)$$

where

$$\begin{aligned}f_1 &\triangleq \sigma \left(2e^2 w |w| + b + \frac{e^2}{2} \right) + q - 1 - \sigma \left(b - \frac{e^2}{2} \right) + p, \\ f_2 &\triangleq \sigma \left(b - \frac{e^2}{2} \right) - p.\end{aligned}$$

From Eq. (69), we see that the second derivatives of L depend on the first-derivatives of f_1 and f_2 , which we now compute. Recall that the derivative of the sigmoid function obeys $\sigma'(z) = \sigma(z)(1 - \sigma(z)) = \sigma(z)\sigma(-z)$ for any $z \in \mathbb{R}$. Now the gradients of f_1 and f_2 with respect to b, e , and w are

$$\begin{aligned}\frac{\partial f_1}{\partial b} &= \sigma \left(2e^2 w |w| + b + \frac{e^2}{2} \right) \sigma \left(-2e^2 w |w| - b - \frac{e^2}{2} \right) - \sigma \left(b - \frac{e^2}{2} \right) \sigma \left(-b + \frac{e^2}{2} \right), \\ \frac{\partial f_2}{\partial b} &= \sigma \left(b - \frac{e^2}{2} \right) \sigma \left(-b + \frac{e^2}{2} \right), \\ \frac{\partial f_1}{\partial e} &= (4ew|w| + e) \sigma \left(2e^2 w |w| + b + \frac{e^2}{2} \right) \sigma \left(-2e^2 w |w| - b - \frac{e^2}{2} \right) + e \sigma \left(b - \frac{e^2}{2} \right) \sigma \left(-b + \frac{e^2}{2} \right), \\ \frac{\partial f_2}{\partial e} &= (-e) \sigma \left(b - \frac{e^2}{2} \right) \sigma \left(-b + \frac{e^2}{2} \right), \\ \frac{\partial f_1}{\partial w} &= (4e^2 \cdot w \text{sign}(w)) \sigma \left(2e^2 w |w| + b + \frac{e^2}{2} \right) \sigma \left(-2e^2 w |w| - b - \frac{e^2}{2} \right), \\ \frac{\partial f_2}{\partial w} &= 0.\end{aligned}$$

Using the fact that $\sigma\left(\log \frac{p}{q}\right) = \frac{p}{p+q} = \pi_1$ and $\sigma\left(-\log \frac{p}{q}\right) = \frac{q}{p+q} = \pi_0$, the above gradients evaluated for any $\gamma = (b = \log \frac{p}{q}, e = 0, w)$ further reduce to

$$\begin{aligned} \frac{\partial f_1}{\partial b} \Big|_{\gamma} &= 0, & \frac{\partial f_2}{\partial b} \Big|_{\gamma} &= \pi_0 \pi_1, \\ \frac{\partial f_1}{\partial e} \Big|_{\gamma} &= 0, & \frac{\partial f_2}{\partial e} \Big|_{\gamma} &= 0, \\ \frac{\partial f_1}{\partial w} \Big|_{\gamma} &= 0, & \frac{\partial f_2}{\partial w} \Big|_{\gamma} &= 0. \end{aligned} \tag{70}$$

Now substituting Eq. (70) when computing the second-derivatives of L in Eq. (69), we obtain

$$\begin{aligned} \frac{\partial^2 L}{\partial b^2} \Big|_{\gamma} &= \mathbb{E}_X \left[\frac{\partial f_1}{\partial b} \Big|_{\gamma} X + \frac{\partial f_2}{\partial b} \Big|_{\gamma} \right] = \pi_0 \pi_1, \\ \frac{\partial^2 L}{\partial b \partial e} \Big|_{\gamma} &= \mathbb{E}_X \left[\frac{\partial f_1}{\partial e} \Big|_{\gamma} X + \frac{\partial f_2}{\partial e} \Big|_{\gamma} \right] = 0, \\ \frac{\partial^2 L}{\partial b \partial w} \Big|_{\gamma} &= \mathbb{E}_X \left[\frac{\partial f_1}{\partial w} \Big|_{\gamma} X + \frac{\partial f_2}{\partial w} \Big|_{\gamma} \right] = 0, \\ \frac{\partial^2 L}{\partial e^2} \Big|_{\gamma} &= \mathbb{E}_X [(f_1 X + f_2)(2X(1 + 2w|w|) - 1)] \Big|_{\gamma} \\ &= \mathbb{E}_X [(2f_1(1 + 2w|w|) - f_1 + 2f_2(1 + 2w|w|)X - f_2)] \Big|_{\gamma} \\ &= \mathbb{E}_X [(f_1(1 + 4w|w|) + f_2(2 + 4w|w|)X - f_2)] \Big|_{\gamma} \\ &= (f_1(1 + 4w|w|) + f_2(2 + 4w|w|))\pi_1 - f_2 \Big|_{\gamma} \\ &\stackrel{(a)}{=} ((p + q - 1)(1 + 4w|w|) - \pi_1(p + q - 1)(2 + 4w|w|))\pi_1 + \pi_1(p + q - 1) \\ &= \pi_1(p + q - 1)(1 + 4w|w| - \pi_1(2 + 4w|w|) + 1) \\ &\stackrel{(b)}{=} 2\pi_1\pi_0(p + q - 1)(1 + 2w|w|), \end{aligned} \tag{71}$$

2035 where (a) follows from the fact that $f_1|_\gamma = p + q - 1$, $f_2|_\gamma = \sigma(b) - p = \frac{p}{p+q} - p = \frac{-p}{p+q}(p + q - 1) = -\pi_1(p + q - 1)$
 2036 and (b) from $1 - \pi_1 = \pi_0$. Returning to the remaining second derivatives,

$$\begin{aligned}
 2037 & \\
 2038 & \frac{\partial^2 L}{\partial e \partial w} \Big|_\gamma = \frac{\partial}{\partial e} (\mathbb{E}[(f_1 X + f_2) X] \cdot 4e^2 |w|) \Big|_\gamma \\
 2039 & = \frac{\partial}{\partial e} (\mathbb{E}[(f_1 + f_2) X] \cdot 4e^2 |w|) \Big|_\gamma \\
 2040 & = \frac{\partial}{\partial e} ((f_1 + f_2) \cdot 4\pi_1 e^2 |w|) \Big|_\gamma \\
 2041 & = \frac{\partial}{\partial e} \left(\left(\sigma \left(2e^2 w |w| + b + \frac{e^2}{2} \right) + q - 1 \right) \cdot 4\pi_1 e^2 |w| \right) \Big|_\gamma \\
 2042 & = \left(\frac{\partial}{\partial e} \sigma \left(2e^2 w |w| + b + \frac{e^2}{2} \right) \right) 4\pi_1 e^2 |w| \Big|_\gamma \\
 2043 & = \left(\frac{\partial}{\partial e} \sigma \left(2e^2 w |w| + b + \frac{e^2}{2} \right) \right) 4\pi_1 e^2 |w| \Big|_\gamma \\
 2044 & \quad + \left(\sigma \left(2e^2 w |w| + b + \frac{e^2}{2} \right) + q - 1 \right) \cdot \frac{\partial}{\partial e} (4\pi_1 e^2 |w|) \Big|_\gamma \\
 2045 & = 0, \\
 2046 & \frac{\partial^2 L}{\partial w^2} \Big|_\gamma = \frac{\partial}{\partial w} (\mathbb{E}[(f_1 X + f_2) X] \cdot 4e^2 |w|) \\
 2047 & = \left(\frac{\partial}{\partial w} \mathbb{E}[(f_1 X + f_2) X] \cdot 4|w| \right) e^2 \Big|_\gamma \\
 2048 & = 0.
 \end{aligned} \tag{72}$$

2049 Congregating all the second derivatives from Eq. (71) and Eq. (72) into the Hessian $\mathbf{H}(\gamma)$ in Eq. (68), we finally obtain

$$\mathbf{H}(\gamma) = \pi_0 \pi_1 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2(p + q - 1)(1 + 2w|w|) & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

□

K. Proofs of lemmas in App. F

K.1. Proof of Lemma 5

Proof. Recall from Lemma 2 that for any $\theta = (e, w) \in \mathbb{R}^2$ and $b \in \mathbb{R}$, we have

$$L(\theta, b) = \mathbb{E}_{X,Y} \left[\ell_{\log} \left((2Y - 1) \left(e^2(1 + 2w|w|)X + b - \frac{e^2}{2} \right) \right) \right].$$

Since $\ell_{\log}(\cdot)$ is a convex function, $Y \in \{0, 1\}$ and thus $2Y - 1 \in \{\pm 1\}$, the convexity of L in b follows from the following fact:

$$\frac{\partial^2 L}{\partial b^2} = \mathbb{E}_{X,Y} \left[\ell''_{\log} \left((2Y - 1) \left(e^2(1 + 2w|w|)X + b - \frac{e^2}{2} \right) \right) \right] \geq 0.$$

To find the optimal b_* , we set the gradient $\frac{\partial L}{\partial b} = 0$. Thus from Lemma 3, we obtain

$$\begin{aligned} \frac{\partial L}{\partial b} &= \mathbb{E}_X [f_1 X + f_2] = 0, \\ f_1 &= \sigma \left(2e^2 w|w| + b + \frac{e^2}{2} \right) + q - 1 - \sigma \left(b - \frac{e^2}{2} \right) + p, f_2 = \sigma \left(b - \frac{e^2}{2} \right) - p. \end{aligned}$$

Substituting $E \triangleq e^2(1 + 2w|w|)$, $B \triangleq b - \frac{e^2}{2}$, and $\mathbb{E}[X] = \pi_1 = p/(p + q)$ in the above equations,

$$\pi_1 (\sigma(E + B) + q - 1 - \sigma(B) + p) + \sigma(B) - p = \pi_1 \cdot \sigma(E + B) + \pi_0 \cdot \sigma(B) + \frac{p(p + q - 1)}{p + q} - p = 0.$$

Further simplifying,

$$\pi_0 \cdot \sigma(B) = \pi_1 \cdot (1 - \sigma(E + B)) = \pi_1 \cdot \sigma(-E - B).$$

In other words,

$$\frac{(1 + \exp(-B))^{-1}}{(1 + \exp(E + B))^{-1}} = \frac{\pi_1}{\pi_0} = \frac{p}{q} \Rightarrow \frac{1 + \exp(E) \exp(B)}{1 + \exp(-B)} = \frac{p}{q}.$$

Defining $x \triangleq \exp(B)$ and $A \triangleq \exp(E)$, we thus obtain the following quadratic equation in x and its corresponding roots:

$$Ax^2 - x \left(\frac{p}{q} - 1 \right) - \frac{p}{q} = 0 \Rightarrow x = \frac{1}{2A} \left[\frac{p}{q} - 1 \pm \sqrt{\left(\frac{p}{q} - 1 \right)^2 + 4 \cdot \frac{p}{q} \cdot A} \right].$$

Since $x > 0$, we take the root corresponding to the addition choice above and resubstituting $x = \exp(b - \frac{e^2}{2})$ and $A = \exp(e^2(1 + 2w|w|))$, we obtain the final expression for b_* . In particular, if $e = 0$, it is easy to see that $A = 1$ and hence $x = \exp(b_*) = \frac{p}{q}$, implying $b_* = \log \frac{p}{q}$. Similarly, if $A = \frac{(1-p)(1-q)}{pq}$, it's straightforward to see that $\exp(b_* - e^2/2) = \frac{p}{1-p}$ and hence $b_* - e^2/2 = \log \frac{p}{1-p}$. \square

K.2. Proof of Lemma 6

Proof. The proof directly follows from Lemma 2 by substituting $b = b_*$. \square

K.3. Proof of Lemma 7

Proof. By Danskin's theorem (Danskin, 1966), it follows that for $b_* = \operatorname{argmin}_{b \in \mathbb{R}} L(\theta, b)$ and $L(\theta) = L(\theta, b_*)$, we have $\nabla_{\theta} L(\theta) = \nabla_{\theta} L(\theta, b_*)$. Using the gradient expressions of $L(\theta, b)$ w.r.t θ from Lemma 3, and using the fact that $\frac{\partial L}{\partial b} = \mathbb{E}[f_1 X + f_2]$ at $b = b_*$, the claim follows. \square

K.4. Proof of Lemma 8

Proof. Since $L(\theta) = L(\theta, b_*)$ where $b_* = \operatorname{argmin}_{b \in \mathbb{R}} L(\theta, b)$, the identity in Eq. (35) about the Hessian of the loss L with respect to θ follows from the classical result of (?)Lemma 2.2]shapiro1985second about second-derivatives of extremal-value functions. Finally Eq. (36) follows from substituting the full Hessian in $\mathbb{R}^{3 \times 3}$ from Lemma 4 in this identity. \square

L. Proofs of lemmas in App. G

L.1. Proof of Lemma 9

Proof. Denote $(e_t, w_t) = (e, w)$ with the dependence on time implicitly assumed. Then by the definition of GF and the gradient expressions in Lemma 7, we have that

$$\frac{de}{dt} = -\frac{\partial L}{\partial e}(\boldsymbol{\theta}_t) = -2\mathbb{E}[(f_1 X + f_2)X] \cdot (1 + 2w|w|)e, \quad (73)$$

$$\frac{dw}{dt} = -\frac{\partial L}{\partial w}(\boldsymbol{\theta}_t) = -\mathbb{E}[(f_1 X + f_2)X] \cdot 4e^2|w|. \quad (74)$$

Dividing Eq. (73) by $(1 + 2w|w|)e$ and Eq. (74) by $4e^2|w|$, we have

$$\frac{1}{1 + 2w|w|} \frac{de}{dt} = \frac{1}{2e|w|} \frac{dw}{dt} \Rightarrow e \frac{de}{dt} = \left(w + \frac{1}{2|w|} \right) \frac{dw}{dt}.$$

Noting that $\frac{d}{dt}(\text{sign}(w) \cdot \log |w|) = \frac{1}{|w|} \frac{dw}{dt}$, the above equation can be rewritten as

$$\frac{d}{dt} (e^2 - w^2 - \text{sign}(w) \cdot \log |w|) = 0.$$

Thus defining $\mathcal{E}(e, w) = e^2 - w^2 - \text{sign}(w) \cdot \log |w|$ for $w \neq 0$, the above equation implies $\mathcal{E}(\boldsymbol{\theta}_t) = \mathcal{E}(\boldsymbol{\theta}_0)$ for $\boldsymbol{\theta}_0 = (e_0, w_0)$ with $w_0 \neq 0$. On the other hand, it's easy to see that if $w_0 = 0$, Eq. (74) implies $\frac{dw}{dt} = 0$ at $t = 0$ and hence $w_t = 0$ for all $t \geq 0$. \square

L.2. Proof of Lemma 10

Proof. To prove the convergence of the trajectory $(\boldsymbol{\theta}_t)_{t \geq 0}$, we use the classical result due to Łojasiewicz (?) Theorem 2.2]absil2005convergence which gurantees the convergence of gradient flow for real analytic functions, as long as the trajectory is bounded. Hence we first show the boundedness of the trajectory.

(i) $(\boldsymbol{\theta}_t)_{t \geq 0}$ is bounded.

We consider the cases $\boldsymbol{\theta}_0 \in e$ -axis and $\boldsymbol{\theta}_0 \in \mathbb{R}^2 \setminus e$ -axis separately.

Let's suppose $\boldsymbol{\theta}_0 = (e_0, w_0) \in e$ -axis, i.e. $w_0 = 0$. Thus it follows from Lemma 9 that $w_t = 0$ for all $t \geq 0$. That is, the trajectory always stays on the e -axis and it suffices to track $(e_t)_{t \geq 0}$ and show that they are bounded. To this end, we show that if $e_0 > 0$, we have $\frac{de}{dt} < 0$ and if $e_0 < 0$, we have $\frac{de}{dt} > 0$ for all $t \geq 0$, which establishes our claim. We have from the GF and Lemma 7 that

$$\frac{de_t}{dt} = -\frac{\partial L}{\partial e}(e_t, w_t = 0) = -2\mathbb{E}_X [(f_1 X + f_2)X] e_t = -2\pi_1(f_1 + f_2)e_t, \quad (75)$$

and

$$\begin{aligned} f_1 + f_2 &= -\frac{f_2}{\pi_1} + f_2 = -\frac{\pi_0}{\pi_1} f_2 = -\frac{\pi_0}{\pi_1} \left[\sigma \left((b_\star)_t - \frac{e_t^2}{2} \right) - p \right] \\ &= -\frac{\pi_0}{\pi_1} \left[\left(1 + \exp \left(-(b_\star)_t + \frac{e_t^2}{2} \right) \right)^{-1} - p \right] \\ &= -\frac{\pi_0}{\pi_1} \left[\left(1 + \frac{2x_t}{\frac{p}{q} - 1 + \sqrt{\left(\frac{p}{q} - 1 \right)^2 + 4 \cdot \frac{p}{q} \cdot x_t}} \right)^{-1} - p \right], \quad x_t \triangleq \exp(e_t^2). \end{aligned} \quad (76)$$

Defining

$$g(x) \triangleq \frac{2x}{\frac{p}{q} - 1 + \sqrt{\left(\frac{p}{q} - 1 \right)^2 + 4 \cdot \frac{p}{q} \cdot x}}, \quad (77)$$

and substituting Eq. (77) and Eq. (76) in Eq. (75), we obtain

$$\frac{de_t}{dt} = 2\pi_0 \left(\frac{1}{1+g(x_t)} - p \right) \cdot e_t. \quad (78)$$

Since $x_t = \exp(e_t^2) = \exp(-e_t^2)$, in view of Eq. (78), with out loss of generality, we can assume that $e_0 > 0$ and show that $\frac{de_t}{dt} < 0$ for all $t \geq 0$. That is, the RHS Eq. (78) is negative. Note that $x_t \geq 1$ since $x_t = \exp(e_t^2)$ and $g(x_t) > 0$ since the denominator $\frac{p}{q} - 1 + \sqrt{\left(\frac{p}{q} - 1\right)^2 + 4 \cdot \frac{p}{q} \cdot x_t} > \frac{p}{q} - 1 + \frac{p}{q} + 1 = 2 \cdot \frac{p}{q} > 0$. Further $g(1) = \frac{q}{p}$ and $\lim_{x \rightarrow \infty} g(x) = \infty$. If we show that $g(x)$ is increasing in x for $x \geq 1$, it implies $\frac{1}{1+g(x)} - p < \frac{1}{1+g(1)} - p = \frac{1}{1+\frac{q}{p}} - p = -\frac{p}{p+q}(p+q-1) < 0$. Thus the gradient in Eq. (78) remains negative starting at $t = 0$ and hence the sequence $(e_t)_{t \geq 0}$ will be bounded. Now it remains to show $g(\cdot)$ is increasing, i.e. $g'(\cdot) > 0$. Defining $C = \left(\frac{p}{q} - 1\right) / \left(2\sqrt{\frac{p}{q}}\right)$ and $D = C^2$, we have that $g(x)$ upto a positive scaling is

$$g(x) = \frac{x}{C + \sqrt{x+D}}.$$

Hence

$$g'(x) = \frac{C + \sqrt{x+D} - \frac{x}{2\sqrt{x+D}}}{(C + \sqrt{x+D})^2}.$$

Thus it suffices to show that $h_1(x) \triangleq C + \sqrt{x+D} > h_2(x) \triangleq \frac{x}{2\sqrt{x+D}}$ for $x \geq 1$. Note that $h_1(1) - h_2(1)$ is given by

$$\begin{aligned} h_1(1) - h_2(1) &= \frac{\frac{p}{q} - 1}{2\sqrt{\frac{p}{q}}} + \sqrt{1 + \left(\frac{\frac{p}{q} - 1}{2\sqrt{\frac{p}{q}}}\right)^2} - \frac{1}{2\sqrt{1 + \left(\frac{\frac{p}{q} - 1}{2\sqrt{\frac{p}{q}}}\right)^2}} \\ &= \sqrt{\frac{p}{q}} - \frac{\sqrt{\frac{p}{q}}}{1 + \frac{p}{q}} \\ &> 0. \end{aligned}$$

Now we show that $h_1'(x) > h_2'(x)$ for all $x \geq 1$ which implies that $h_1(x) > h_2(x)$ for all $x \geq 1$, thus establishing our claim. To this end, we have that

$$\begin{aligned} h_1'(x) - h_2'(x) &= \frac{1}{2\sqrt{x+D}} - \frac{\sqrt{x+D} - \frac{x}{2\sqrt{x+D}}}{2(x+D)} \\ &= \frac{x}{2\sqrt{x+D}(x+D)} > 0. \end{aligned}$$

This proves our claim that $g(\cdot)$ is increasing and hence $(e_t)_{t \geq 0}$, and consequently $(\theta_t)_{t \geq 0}$, is bounded when $\theta_0 \in e$ -axis.

Now let's assume that $\theta_0 = (e_0, w_0) \in \mathbb{R}^2 \setminus e$ -axis. Since $(\theta_t)_{t \geq 0} \subseteq \mathbb{R}^2 \setminus e$ -axis, it follows that the loss $L(\cdot)$ is analytic on the trajectory (since the logistic function is analytic), and hence by (?)Theorem 2.2]absil2005convergence, it follows that $\lim_{t \rightarrow \infty} \|\theta_t\|$ exists. Now we show that $\lim_{t \rightarrow \infty} \|\theta_t\| \neq \infty$, which implies the desired result about boundedness. To show $\lim_{t \rightarrow \infty} \|\theta_t\| \neq \infty$, we show that there exists a large $B > 0$ such that for any $\theta_t = (e, w) \in \mathbb{R}^2$ with $\|(e, w)\| \geq B$, the velocity vector $\frac{d\theta}{dt}$ points inwards into the ball of radius B and thus the trajectory always stays inside this ball, and hence bounded. To establish this, let's denote $(e_t, w_t) = (e, w)$ with the dependence on time implicitly assumed. Then by the definition of GF and the gradient expressions in Lemma 7, we have that

$$\frac{de}{dt} = -\frac{\partial L}{\partial e}(\theta_t) = -2\mathbb{E}[(f_1 X + f_2)X] \cdot (1 + 2w|w|)e \quad (79)$$

$$\frac{dw}{dt} = -\frac{\partial L}{\partial w}(\theta_t) = -\mathbb{E}[(f_1 X + f_2)X] \cdot 4e^2|w|, \quad (80)$$

where $f_1 = \sigma\left(2e^2w|w| + b_* + \frac{e^2}{2}\right) + q - 1 - \sigma\left(b_* - \frac{e^2}{2}\right) + p$, and $f_2 = \sigma\left(b_* - \frac{e^2}{2}\right) - p$ with $\pi_1 f_1 + f_2 = 0$. Given that only $\frac{de}{dt}$ flips in sign under the transformation $(e, w) \mapsto (-e, w)$, with out loss of generality we can assume $e > 0$. Now let's also assume $w > 0$. Thus, in view of Eq. (79) and GF, to show that the derivative points inwards, it suffices to show that $\mathbb{E}[(f_1 X + f_2)X] > 0$ for reasonably large B with $\|(e, w)\| = B$. Similar to Eq. (76) and Eq. (77) above, using the relation $\pi_1 f_1 + f_2 = 0$, we obtain

$$\mathbb{E}[(f_1 X + f_2)X] = \pi_1(f_1 + f_2) = -\pi_0 \left(\frac{1}{1+g(x)} - p \right), \quad x \triangleq \exp(e^2(1+2w|w|)). \quad (81)$$

Using the fact that $g(x)$ is increasing for $x \geq 1$ with $\lim_{x \rightarrow \infty} g(x) = \infty$, and $|w| = w > 0$, we can chose a $B > 0$ such that for any $\|(e, w)\| \geq B$, in Eq. (81) we have $1/(1+g(x)) < p$ and hence $\mathbb{E}[(f_1 X + f_2)X] > 0$. This finishes the proof of our claim. The proof for $w < 0$ is similar, where we make use of the fact that $\lim_{x \rightarrow 0} g(x) = 0$ to show $\mathbb{E}[(f_1 X + f_2)X] < 0$ for e, w reasonably large.

(ii) $\lim_{t \rightarrow \infty} \theta_t = \theta_{\text{lim}}$. Since the logistic function $\ell_{\log}(\cdot)$ is analytic, it follows from Lemma 6 that the loss $L(\theta)$ is analytic too whenever $w \neq 0$. On the other hand, when $w = 0$, it's easy to see that L is an analytic function of $e \in \mathbb{R}$. By Lemma 9, we know that if $w_0 \neq 0$, $w_t \neq 0$ and if $w_0 = 0$, $w_t = 0$ for all $t \geq 0$. Thus the loss is analytic on the trajectory for all $t \geq 0$. Since the trajectory is bounded, it follows from Łojasiewicz's theorem (Theorem 2.2) that there exists a $\theta_{\text{lim}} \in \mathbb{R}^2$ such that $\lim_{t \rightarrow \infty} \theta_t = \theta_{\text{lim}}$.

(iii) $\lim_{t \rightarrow \infty} \|\nabla L(\theta_t)\| = \|\nabla L(\theta_{\text{lim}})\| = 0$. Since the trajectory is bounded, it follows from Theorem 2 (Ahmadova 2023) convergence that the gradient converges to zero, i.e. $\lim_{t \rightarrow \infty} \|\nabla L(\theta_t)\| = 0$. Since $\nabla L(\cdot)$ is a continuous function and $\lim_{t \rightarrow \infty} \theta_t = \theta_{\text{lim}}$, we have $\lim_{t \rightarrow \infty} \|\nabla L(\theta_t)\| = \|\nabla L(\theta_{\text{lim}})\| = 0$. □

L.3. Proof of Lemma 11

Proof. Since the energy function $\mathcal{E}(\cdot, \cdot)$ in Eq. (51) is a continuous function in $\mathbb{R}^2 \setminus e\text{-axis}$, and any trajectory $(\theta_t)_{t \geq 0}$ with initialization $\theta \in \mathbb{R}^2 \setminus e\text{-axis}$ stays in $\mathbb{R}^2 \setminus e\text{-axis}$ for all $t \geq 0$ (Lemma 9), it follows that $\lim_{t \rightarrow \infty} \mathcal{E}(\theta_t) = \mathcal{E}(\theta_{\text{lim}}) = \mathcal{E}(\theta_0)$. As $\nabla L(\theta_{\text{lim}}) = 0$ from Lemma 10, it follows that θ_{lim} lies at the intersection of the contour line $\mathcal{E}(e, w) = \mathcal{E}_0$ with the set of critical points of L in \mathbb{R}^2 .

On the other hand, if $\theta_0 \in e\text{-axis}$, we have $\theta_t \in e\text{-axis}$ from Lemma 9 for all $t \geq 0$. Hence $\theta_{\text{lim}} \in e\text{-axis}$. □

L.4. Proof of Lemma 12

Proof. Recall that $f : \mathbb{R} \setminus \{0\} \rightarrow \mathbb{R}$, defined as $f(w) \triangleq \mathcal{E}(e=0, w) = -(w^2 + \text{sign}(w) \cdot \log|w|)$. If $w < 0$, we have

$$f(w) = -(w^2 - \log(-w)), \quad f'(w) = -2w + \frac{1}{w}.$$

Hence $f'(w) \geq 0$ for $w \in (-\infty, -1/\sqrt{2}]$ and $f'(w) \leq 0$ for $w \in [-1/\sqrt{2}, 0)$ with $f'(-1/\sqrt{2}) = 0$. It's also straightforward to see that $\lim_{w \rightarrow -\infty} f(w) = -\infty$, $\lim_{w \rightarrow 0^-} f(w) = -\infty$, and $f(-1/\sqrt{2}) = \mathcal{E}_{\text{sad}}$ (by the definition of f). This establishes (i), (ii), and (iii).

On the other hand, for $w > 0$, we have $f(w) = -(w^2 + \log w)$ and $f'(w) = -(2w + 1/w)$. Hence f is monotonically decreasing for $w > 0$ with $\lim_{w \rightarrow 0^+} f(w) = \infty$ and $\lim_{w \rightarrow \infty} f(w) = -\infty$. Note that $w = 0$ acts as an energy barrier since $\lim_{w \rightarrow 0^-} f(w) = -\infty$ whereas $\lim_{w \rightarrow 0^+} f(w) = \infty$. □

M. Proofs of lemmas in App. H
M.1. Proof of Lemma 13

Proof. First we consider the loss with the bias $L(\boldsymbol{\theta}, b)$ from 44:

$$L(\boldsymbol{\theta}, b) = \mathbb{E}_{X,Y} [\ell_{\log}((2Y - 1) \cdot \text{logit}_X(\boldsymbol{\theta}, b))] \\ \text{logit}_X(\boldsymbol{\theta}, b) = e^2 \left[\left(X - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] + b, \quad \boldsymbol{\theta} = (e, w, a).$$

Using the fact that $\ell'_{\log}(z) = \sigma(z) - 1$ and $2Y - 1 \in \{\pm 1\}$, we have for any $\theta \in \{e, w, a, b\}$ that

$$\nabla_{\theta} L = \mathbb{E}[(\sigma((2Y - 1) \cdot \text{logit}_X) - 1)(2Y - 1) \cdot \nabla_{\theta} \text{logit}_X] = \mathbb{E}[(\sigma(\text{logit}_X) - Y) \cdot \nabla_{\theta} \text{logit}_X]. \quad (82)$$

Now we simplify $\sigma(\text{logit}_X)$ using the fact that $X \in \{0, 1\}$:

$$\begin{aligned} \sigma(\text{logit}_X) &= \sigma \left(e^2 \left[\left(X - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] + b \right) \\ &= X \underbrace{\sigma \left(e^2 \left[\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] + b \right)}_{\triangleq \phi_1} \\ &\quad + (1 - X) \underbrace{\sigma \left(e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] + b \right)}_{\triangleq \phi_0} \\ &= X\phi_1 + (1 - X)\phi_0 \\ &= X(\phi_1 - \phi_0) + \phi_0. \end{aligned}$$

Thus the gradients in 82 are given by

$$\begin{aligned} \nabla_{\theta} L &= -\mathbb{E}[(Y - X(\phi_1 - \phi_0) - \phi_0) \nabla_{\theta} \text{logit}_X] \\ &= -\mathbb{E}_X \left[\mathbb{E}_{x_1^{n+1}} [(\mathbb{E}[Y | X] - X(\phi_1 - \phi_0) - \phi_0) \nabla_{\theta} \text{logit}_X] \right] \\ &= -\mathbb{E}_X [(1 - p - q)X + p - X(\phi_1 - \phi_0) - \phi_0] \nabla_{\theta} \text{logit}_X \\ &= -\mathbb{E}_X \left[\left(\underbrace{\left((1 - p - q - \phi_1 + \phi_0) \right)}_{f_1} X + \underbrace{p - \phi_0}_{f_2} \right) \nabla_{\theta} \text{logit}_X \right] \\ &= -\mathbb{E}_X [(f_1 X + f_2) \nabla_{\theta} \text{logit}_X]. \end{aligned} \quad (83)$$

Now we compute the individual gradients with respect to e, w, a and b . Recall that

$$\text{logit}_X = e^2 \left[\left(X - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] + b.$$

2365 Thus,

$$\begin{aligned}
 2366 \quad \nabla_e \text{logit}_X &= 2e \left[\left(X - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] \\
 2367 &+ e^2 \left[2ae \left(X - \frac{1}{2} \right) (1 + 2w|w|) + w \text{sign}(w(1 + ae^2)) (2ae) \right] \\
 2368 &= 2e \left(X - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + 2ew|w|(1 + ae^2) \\
 2369 &+ 2e^3 a \left(X - \frac{1}{2} \right) (1 + 2w|w|) + 2e^3 aw \text{sign}(w(1 + ae^2)) \\
 2370 &= 2e \left(X - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + 2e^3 a \left(X - \frac{1}{2} \right) (1 + 2w|w|) \\
 2371 &+ 2ew|w|(1 + ae^2) + 2e^3 aw \text{sign}(w(1 + ae^2)).
 \end{aligned}$$

2372 Substituting the above equation in Eq. (83), we obtain

$$\begin{aligned}
 2373 \quad \nabla_e L &= - \left(\mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \right) \cdot 2e (1 + ae^2) (1 + 2w|w|) \\
 2374 &- \left(\mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \right) \cdot 2e^3 a (1 + 2w|w|) \\
 2375 &- (\mathbb{E} [(f_1 X + f_2)]) \cdot 2ew|w|(1 + ae^2) \\
 2376 &- (\mathbb{E} [(f_1 X + f_2)]) \cdot 2e^3 aw \text{sign}(w(1 + ae^2)).
 \end{aligned} \tag{84}$$

2377 Now we compute the derivative with respect to w .

$$\begin{aligned}
 2378 \quad \text{logit}_X &= e^2 \left[\left(X - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] + b \\
 2379 \Rightarrow \nabla_w \text{logit}_X &= 2e^2 \left(X - \frac{1}{2} \right) (1 + ae^2) (|w| + \text{sign}(w) w) \\
 2380 &+ e^2 [w(1 + ae^2) + w(1 + ae^2) \text{sign}(w(1 + ae^2))] \\
 2381 \Rightarrow \nabla_w L &= - \left(\mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \right) 2e^2 (1 + ae^2) (|w| + \text{sign}(w) w) \\
 2382 &- (\mathbb{E} [(f_1 X + f_2)]) e^2 [w(1 + ae^2) + w(1 + ae^2) \text{sign}(w(1 + ae^2))].
 \end{aligned}$$

2383 Similarly, for a :

$$\begin{aligned}
 2384 \quad \text{logit}_X &= e^2 \left[\left(X - \frac{1}{2} \right) (1 + 2w|w|) + w|w|(1 + ae^2) \right] + b \\
 2385 \Rightarrow \nabla_a \text{logit}_X &= e^4 \left(X - \frac{1}{2} \right) (1 + ae^2) (1 + 2w|w|) \\
 2386 &+ e^4 w^2 \text{sign}(w(1 + ae^2)) \\
 2387 \Rightarrow \nabla_a L &= - \left(\mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \right) e^4 (1 + 2w|w|) \\
 2388 &- (\mathbb{E} [(f_1 X + f_2)]) e^4 w^2 \text{sign}(w(1 + ae^2)).
 \end{aligned}$$

2389 Finally, since $\nabla_b \text{logit}_X = 1$, it follows from Eq. (83) that

$$2390 \quad \nabla_b L = -\mathbb{E} [f_1 X + f_2]. \tag{85}$$

2391 For the optimal b_* , we have $\nabla_b L = 0$ and hence $\mathbb{E} [f_1 X + f_2] = 0$, simplifying the expressions for the gradients of e , w , and a above. This concludes the proof.

2392 □

2420 **M.2. Optimality conditions for linear self-attention**

 2421 We prove Thm. 9 and Thm. 10 below. Note that Thm. 10 directly follows from the former by removing the bias b since for
 2422 critical points, the bias b is already the optimal one, similar to the proof of Thm. 1.
 2423

2424 We characterize the set of global minima first.

 2425 **Set of all global minima.** Let $\gamma_\star \in \mathbb{R}^4$ be arbitrary. From (Makkuva et al., 2024)-Lemma 1, we have that γ_\star is a global
 2426 minimum for the loss $L(\cdot)$ in Eq. (26) if and only if its prediction probability satisfies $f_{\gamma_\star}(x_1^n) = \mathbb{P}(x_{n+1} = 1 | x_n)$, the
 2427 Markov kernel. Since the input $\{x_n\}_{n=1}^N \sim (\pi(p, q), \mathbf{P}(p, q))$, we have that
 2428

$$2429 \mathbb{P}(x_{n+1} = 1 | x_n) = (1 - x_n)p + x_n(1 - q) = (1 - p - q)x_n + p. \quad (86)$$

 2431 On the other hand, by definition, from Eq. (7), $f_{\gamma_\star}(x_1^n) = \sigma\left(e^2 \left[\left(x_n - \frac{1}{2}\right) (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right)$,
 2432 where $\gamma_\star = (e, w, b)$. Since $x_n \in \{0, 1\}$, this can be further simplified to
 2433

$$2434 f_{\gamma_\star}(x_1^n) = \sigma\left(e^2 \left[\left(x_n - \frac{1}{2}\right) (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right) \quad (87)$$

$$2435 = x_n \sigma\left(e^2 \left[\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right) \quad (88)$$

$$2436 + (1 - x_n) \sigma\left(e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right)$$

$$2437 = x_n \sigma\left(e^2 \left[\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right)$$

$$2438 - x_n \sigma\left(e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right) \quad (89)$$

$$2439 + \sigma\left(e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right)$$

 2440 Since both $f_{\gamma_\star}(x_1^n)$ and $\mathbb{P}(x_{n+1} = 1 | x_n)$ are linear functions of x_n , equating them for all values of $x_n \in \{0, 1\}$ implies
 2441 that the respective coefficients in these functions in Eq. (86) and Eq. (89) are also equal, i.e.
 2442

$$2443 1 - p - q = \sigma\left(e^2 \left[\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right)$$

$$2444 - \sigma\left(e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right)$$

$$2445 p = \sigma\left(e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right)$$

2446 and hence

$$2447 \sigma\left(e^2 \left[\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right) = 1 - q \quad (90)$$

$$2448 \sigma\left(e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b\right) = p$$

 2449 Since $\sigma(z) = y$ for $y \in (0, 1)$ implies $z = \log \frac{y}{1-y}$, Eq. (90) can be rewritten as
 2450

$$2451 e^2 \left[\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)|\right] + b = \log \frac{1 - q}{q}$$

$$e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] + b = \log \frac{p}{1-p}$$

Adding and subtracting the above two equations, we obtain:

$$\begin{aligned} e^2 w|w|(1 + ae^2) + b &= \frac{1}{2} \log \frac{p(1-q)}{q(1-p)} \\ e^2 (1 + ae^2) (1 + 2w|w|) &= \log \frac{(1-q)(1-p)}{pq} \end{aligned} \quad (91)$$

Thus $\gamma_* \in \mathbb{R}^3$ is a global minimum for $L(\cdot)$ if and only if it satisfies Eq. (91) (note that it's already a critical point, as established in Thm. 6). Thus, the set of all global minimum $\Gamma_*(p, q)$ is given by

$$\begin{aligned} \Gamma_*(p, q) \triangleq \{ \gamma_* = (e, w, b, a) \in \mathbb{R}^4 : e^2 w|w|(1 + ae^2) + b &= \frac{1}{2} \log \frac{p(1-q)}{q(1-p)}, \\ e^2 (1 + ae^2) (1 + 2w|w|) &= \log \frac{(1-q)(1-p)}{pq} \} \end{aligned}$$

Since the prediction $f_{\gamma_*}(\cdot)$ equals the Markov kernel for any $\gamma_* \in \Gamma_*$, it follows from Thm. 4 (or (?)Lemma 1]makkuva2024attention) that $L(\gamma_*) = H(x_{n+1} | x_n)$, the entropy rate of the Markov chain.

M.3. Stationary points

Reproducing the derivatives from the previous equations (and setting them to zero)

$$\begin{aligned} \frac{L(\theta)}{\partial b} &= \mathbb{E}_X [f_1 X + f_2] = 0 \\ \frac{L(\theta)}{\partial e} &= \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 2e (1 + ae^2) (1 + 2w|w|) \\ &\quad + \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 2e^3 a (1 + 2w|w|) = 0 \\ \frac{L(\theta)}{\partial w} &= \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 2e^2 (1 + ae^2) (|w| + \text{sign}(w) w) = 0 \\ \frac{L(\theta)}{\partial a} &= \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] e^4 (1 + 2w|w|) = 0 \end{aligned} \quad (92)$$

From the above equations, there are multiple regions of stationarity. Here are the following regions of stationarity:

1. $\mathbb{E}_X [f_1 X + f_2] = 0, e = 0$
2. $\mathbb{E}_X [f_1 X + f_2] = 0, e \neq 0, 1 + ae^2 = 0, 1 + 2w|w| = 0$

Slightly changing the variable order, for any $\gamma = (b, e, w, a) \in \mathbb{R}^4$, we define

$$\mathbf{H}(\gamma) \triangleq \nabla^2 L(\gamma) = \begin{bmatrix} \frac{\partial^2 L}{\partial b^2} & \frac{\partial^2 L}{\partial b \partial e} & \frac{\partial^2 L}{\partial b \partial w} & \frac{\partial^2 L}{\partial b \partial a} \\ \frac{\partial^2 L}{\partial e \partial b} & \frac{\partial^2 L}{\partial e^2} & \frac{\partial^2 L}{\partial e \partial w} & \frac{\partial^2 L}{\partial e \partial a} \\ \frac{\partial^2 L}{\partial w \partial b} & \frac{\partial^2 L}{\partial w \partial e} & \frac{\partial^2 L}{\partial w^2} & \frac{\partial^2 L}{\partial w \partial a} \\ \frac{\partial^2 L}{\partial a \partial b} & \frac{\partial^2 L}{\partial a \partial e} & \frac{\partial^2 L}{\partial a \partial w} & \frac{\partial^2 L}{\partial a^2} \end{bmatrix} \in \mathbb{R}^{4 \times 4}. \quad (93)$$

2530 Note that:

2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584

$$\begin{aligned}
 f_1 &= \sigma \left(\underbrace{e^2 \left[\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)| \right]}_{z_1} + b \right) \\
 &\quad - \sigma \left(\underbrace{e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)| \right]}_{z_2} + b \right) \\
 &\quad + p + q - 1 \\
 f_2 &= \sigma \left(\underbrace{e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w(1 + ae^2)| \right]}_{z_2} + b \right) - p
 \end{aligned}$$

From Eq. (92), we see that the second derivatives of L depend on the first-derivatives of f_1 and f_2 , which we now compute. Recall that the derivative of the sigmoid function obeys $\sigma'(z) = \sigma(z)(1 - \sigma(z)) = \sigma(z)\sigma(-z)$ for any $z \in \mathbb{R}$. Now the gradients of f_1 and f_2 with respect to b, e, w and a are

$$\begin{aligned}
 2585 & \\
 2586 & \\
 2587 & \frac{\partial f_1}{\partial b} = \sigma(z_1)\sigma(-z_1) - \sigma(z_2)\sigma(-z_2) \\
 2588 & \\
 2589 & \frac{\partial f_2}{\partial b} = \sigma(z_2)\sigma(-z_2) \\
 2590 & \\
 2591 & \frac{\partial f_1}{\partial e} = \sigma(z_1)\sigma(-z_1) \left\{ 2e \left[\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] \right\} \\
 2592 & \\
 2593 & \quad + \sigma(z_1)\sigma(-z_1) \left\{ 2ae^3 \left[\frac{1}{2} (1 + 2w|w|) + w|w|\text{sign}(1 + 2w|w|) \right] \right\} \\
 2594 & \\
 2595 & \quad - \sigma(z_2)\sigma(-z_2) \left\{ 2e \left[-\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] \right\} \\
 2596 & \\
 2597 & \quad - \sigma(z_2)\sigma(-z_2) \left\{ 2ae^3 \left[-\frac{1}{2} (1 + 2w|w|) + w|w|\text{sign}(1 + 2w|w|) \right] \right\} \\
 2598 & \\
 2599 & \\
 2600 & \frac{\partial f_2}{\partial e} = \sigma(z_2)\sigma(-z_2) \left\{ 2e \left[-\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right] \right\} \\
 2601 & \\
 2602 & \quad + \sigma(z_2)\sigma(-z_2) \left\{ 2ae^3 \left[-\frac{1}{2} (1 + 2w|w|) + w|w|\text{sign}(1 + 2w|w|) \right] \right\} \\
 2603 & \tag{94} \\
 2604 & \\
 2605 & \frac{\partial f_1}{\partial w} = \sigma(z_1)\sigma(-z_1) \left\{ 2e^2 \left[\frac{1}{2} (1 + ae^2) |w| + |w||1 + ae^2| \right] \right\} \\
 2606 & \\
 2607 & \quad - \sigma(z_2)\sigma(-z_2) \left\{ 2e^2 \left[-\frac{1}{2} (1 + ae^2) |w| + |w||1 + ae^2| \right] \right\} \\
 2608 & \\
 2609 & \\
 2610 & \frac{\partial f_2}{\partial w} = \left\{ 2e^2 \left[-\frac{1}{2} (1 + ae^2) |w| + |w||1 + ae^2| \right] \right\} \\
 2611 & \\
 2612 & \frac{\partial f_1}{\partial a} = \sigma(z_1)\sigma(-z_1) \left\{ e^4 \left[\frac{1}{2} (1 + 2w|w|) + w|w|\text{sign}(1 + ae^2) \right] \right\} \\
 2613 & \\
 2614 & \quad - \sigma(z_2)\sigma(-z_2) \left\{ e^4 \left[-\frac{1}{2} (1 + 2w|w|) + w|w|\text{sign}(1 + ae^2) \right] \right\} \\
 2615 & \\
 2616 & \\
 2617 & \frac{\partial f_2}{\partial a} = \sigma(z_2)\sigma(-z_2) \left\{ e^4 \left[-\frac{1}{2} (1 + 2w|w|) + w|w|\text{sign}(1 + ae^2) \right] \right\} \\
 2618 & \\
 2619 & \\
 2620 &
 \end{aligned}$$

M.3.1. STATIONARY POINTS WHERE $\mathbb{E}_X [f_1 X + f_2] = 0, e = 0$

When $e = 0 \implies z_1 = b, z_2 = b$. Hence:

$$\begin{aligned}
 2624 & f_1 = \sigma(b) + p + q - 1 - \sigma(b) = p + q - 1 \\
 2625 & f_2 = \sigma(b) - p
 \end{aligned}$$

On computing the expectation:

$$\begin{aligned}
 2631 & \mathbb{E}_X [f_1 X + f_2] = (p + q - 1) \mathbb{E}_X [X] + \sigma(b) - p = 0 \\
 2632 & = (p + q - 1) \pi_1 + \sigma(b) - p = 0
 \end{aligned}$$

Rearranging and simplifying:

$$\sigma(b) = \frac{p}{p + q}$$

$$\implies b = \log \frac{p}{q}$$

Using the fact that $\sigma\left(\log \frac{p}{q}\right) = \frac{p}{p+q} = \pi_1$ and $\sigma\left(-\log \frac{p}{q}\right) = \frac{q}{p+q} = \pi_0$, the above gradients evaluated for any $\gamma = (b = \log \frac{p}{q}, e = 0, w, a)$ further reduce to

$$\begin{aligned} \frac{\partial f_1}{\partial b} \Big|_{\gamma} &= 0, & \frac{\partial f_2}{\partial b} \Big|_{\gamma} &= \pi_0 \pi_1, \\ \frac{\partial f_1}{\partial e} \Big|_{\gamma} &= 0, & \frac{\partial f_2}{\partial e} \Big|_{\gamma} &= 0, \\ \frac{\partial f_1}{\partial w} \Big|_{\gamma} &= 0, & \frac{\partial f_2}{\partial w} \Big|_{\gamma} &= 0, \\ \frac{\partial f_1}{\partial a} \Big|_{\gamma} &= 0, & \frac{\partial f_2}{\partial a} \Big|_{\gamma} &= 0. \end{aligned} \tag{95}$$

Now substituting Eq. (95) when computing the second-derivatives of L in Eq. (92), we obtain

$$\begin{aligned} \frac{\partial^2 L}{\partial b^2} \Big|_{\gamma} &= \mathbb{E}_X \left[\frac{\partial f_1}{\partial b} \Big|_{\gamma} X + \frac{\partial f_2}{\partial b} \Big|_{\gamma} \right] = \pi_0 \pi_1, \\ \frac{\partial^2 L}{\partial b \partial e} \Big|_{\gamma} &= \mathbb{E}_X \left[\frac{\partial f_1}{\partial e} \Big|_{\gamma} X + \frac{\partial f_2}{\partial e} \Big|_{\gamma} \right] = 0, \\ \frac{\partial^2 L}{\partial b \partial w} \Big|_{\gamma} &= \mathbb{E}_X \left[\frac{\partial f_1}{\partial w} \Big|_{\gamma} X + \frac{\partial f_2}{\partial w} \Big|_{\gamma} \right] = 0, \\ \frac{\partial^2 L}{\partial b \partial a} \Big|_{\gamma} &= \mathbb{E}_X \left[\frac{\partial f_1}{\partial a} \Big|_{\gamma} X + \frac{\partial f_2}{\partial a} \Big|_{\gamma} \right] \\ \frac{\partial^2 L}{\partial e^2} \Big|_{\gamma} &= \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 2(1 + ae^2)(1 + 2w|w|) \Big|_{\gamma} \\ &\quad + \underbrace{\mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 4e^2 a (1 + 2w|w|)}_0 \Big|_{\gamma} \\ &\quad + \underbrace{\mathbb{E}_X \left[\left(\frac{\partial f_1}{\partial e} \Big|_{\gamma} X + \frac{\partial f_2}{\partial e} \Big|_{\gamma} \right) \left(X - \frac{1}{2} \right) \right] 2(1 + ae^2)(1 + 2w|w|)}_0 \\ &= \mathbb{E}_X \left[(2f_1(1 + 2w|w|) - f_1 + 2f_2(1 + 2w|w|)X - f_2) \right] \Big|_{\gamma} \\ &= \mathbb{E}_X \left[(f_1(1 + 4w|w|) + f_2(2 + 4w|w|)X - f_2) \right] \Big|_{\gamma} \\ &= (f_1(1 + 4w|w|) + f_2(2 + 4w|w|))\pi_1 - f_2 \Big|_{\gamma} \\ &\stackrel{(a)}{=} ((p + q - 1)(1 + 4w|w|) - \pi_1(p + q - 1)(2 + 4w|w|))\pi_1 + \pi_1(p + q - 1) \\ &= \pi_1(p + q - 1)(1 + 4w|w| - \pi_1(2 + 4w|w|) + 1) \\ &\stackrel{(b)}{=} 2\pi_1\pi_0(p + q - 1)(1 + 2w|w|), \end{aligned} \tag{96}$$

2695 where (a) follows from the fact that $f_1|_\gamma = p + q - 1$, $f_2|_\gamma = \sigma(b) - p = \frac{p}{p+q} - p = \frac{-p}{p+q}(p + q - 1) = -\pi_1(p + q - 1)$
 2696 and (b) from $1 - \pi_1 = \pi_0$. Returning to the remaining second derivatives,
 2697
 2698
 2699
 2700
 2701
 2702

$$\begin{aligned}
 2703 \quad \frac{\partial^2 L}{\partial e \partial w} \Big|_\gamma &= \frac{\partial}{\partial e} \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 4e^2 (1 + ae^2) |w| \Big|_\gamma \\
 2704 &= \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] (8e (1 + ae^2) |w| + 8e^3 a |w|) \Big|_\gamma \\
 2705 &+ \mathbb{E}_X \left[\left(\frac{\partial f_1}{\partial e} \Big|_\gamma X + \frac{\partial f_2}{\partial e} \Big|_\gamma \right) \left(X - \frac{1}{2} \right) \right] 4e^2 (1 + ae^2) (|w|) \\
 2706 &= 0 \\
 2707 \quad \frac{\partial^2 L}{\partial e \partial a} \Big|_\gamma &= \frac{\partial}{\partial e} \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] e^4 (1 + 2w|w|) \Big|_\gamma \\
 2708 &= \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 4e^3 (1 + 2w|w|) \Big|_\gamma \\
 2709 &+ \mathbb{E}_X \left[\left(\frac{\partial f_1}{\partial e} \Big|_\gamma X + \frac{\partial f_2}{\partial e} \Big|_\gamma \right) \left(X - \frac{1}{2} \right) \right] e^4 (1 + 2w|w|) \\
 2710 &= 0 \\
 2711 \quad \frac{\partial^2 L}{\partial w^2} \Big|_\gamma &= \frac{\partial}{\partial w} \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 4e^2 (1 + ae^2) |w| \Big|_\gamma \\
 2712 &= \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 4e^2 (1 + ae^2) \text{sign}(w) \\
 2713 &+ \mathbb{E}_X \left[\left(\frac{\partial f_1}{\partial w} \Big|_\gamma X + \frac{\partial f_2}{\partial w} \Big|_\gamma \right) \left(X - \frac{1}{2} \right) \right] 4e^2 (1 + ae^2) |w| \\
 2714 &= 0 \\
 2715 \quad \frac{\partial^2 L}{\partial w \partial a} \Big|_\gamma &= \frac{\partial}{\partial w} \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] e^4 (1 + 2w|w|) \Big|_\gamma \\
 2716 &= \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 2e^4 |w| \Big|_\gamma \\
 2717 &+ \mathbb{E}_X \left[\left(\frac{\partial f_1}{\partial w} \Big|_\gamma X + \frac{\partial f_2}{\partial w} \Big|_\gamma \right) \right] e^4 (1 + 2w|w|) \\
 2718 &= 0 \\
 2719 \quad \frac{\partial^2 L}{\partial a^2} \Big|_\gamma &= \frac{\partial}{\partial a} \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] e^4 (1 + 2w|w|) \Big|_\gamma \\
 2720 &= \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 2e^4 |w| \Big|_\gamma \\
 2721 &+ \mathbb{E}_X \left[\left(\frac{\partial f_1}{\partial a} \Big|_\gamma X + \frac{\partial f_2}{\partial a} \Big|_\gamma \right) \right] e^4 (1 + 2w|w|) \\
 2722 &= 0 \\
 2723 &= 0 \\
 2724 &= 0 \\
 2725 &= 0 \\
 2726 &= 0 \\
 2727 &= 0 \\
 2728 &= 0 \\
 2729 &= 0 \\
 2730 &= 0 \\
 2731 &= 0 \\
 2732 &= 0 \\
 2733 &= 0 \\
 2734 &= 0 \\
 2735 &= 0 \\
 2736 &= 0 \\
 2737 &= 0 \\
 2738 &= 0 \\
 2739 &= 0 \\
 2740 &= 0 \\
 2741 &= 0 \\
 2742 &= 0 \\
 2743 &= 0 \\
 2744 &= 0 \\
 2745 &= 0 \\
 2746 &= 0 \\
 2747 &= 0 \\
 2748 &= 0 \\
 2749 &= 0
 \end{aligned} \tag{97}$$

2750 Congregating all the second derivatives from Eq. (96) and Eq. (97) into the Hessian $\mathbf{H}(\gamma)$ in Eq. (93), we finally obtain

$$2751 \quad 2752 \quad 2753 \quad 2754 \quad 2755 \quad 2756 \quad \mathbf{H}(\gamma) = \pi_0 \pi_1 \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 2(p+q-1)(1+2w|w|) & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

2757 It is trivial to see that (borrowing the result proved in E)

2758 Similar to the notation defined before $\Gamma_{\min}(p, q) \subseteq \mathbb{R}^4$ and $\Gamma_{\text{sad}} \subseteq \mathbb{R}^4$ as follows:

$$2760 \quad 2761 \quad 2762 \quad 2763 \quad 2764 \quad 2765 \quad \Gamma_{\min}(p, q) \triangleq \left\{ \gamma_{\min} = (e, w, b, a) \in \mathbb{R}^4 : e = 0, (p+q-1)(1+2w|w|) > 0, b = \log \frac{p}{q} \right\},$$

$$2766 \quad 2767 \quad 2768 \quad 2769 \quad 2770 \quad \Gamma_{\text{sad}}(p, q) \triangleq \left\{ \gamma_{\text{sad}} = (e, w, b, a) \in \mathbb{R}^4 : e = 0, (p+q-1)(1+2w|w|) \leq 0, b = \log \frac{p}{q} \right\}.$$

2766 M.3.2. STATIONARY POINTS WHERE $\mathbb{E}_X [f_1 X + f_2] = 0, e \neq 0, 1 + ae^2 = 0, 1 + 2w|w| = 0$

2768 For this set of points, the Hessian remains undefined because $\frac{\partial f_1}{\partial e}, \frac{\partial f_2}{\partial e}, \frac{\partial f_1}{\partial a}, \frac{\partial f_2}{\partial a}$ do not exist. This non-existence arises since
2769 sign $(1 + ae^2)$ lacks definition when $1 + ae^2 = 0$. However, even in this scenario, when $e \neq 0, 1 + ae^2 = 0, 1 + 2w|w| =$
2770 $0 \implies z_1 = b, z_2 = b$. Hence:

$$2771 \quad 2772 \quad 2773 \quad 2774 \quad 2775 \quad 2776 \quad f_1 = \sigma(b) + p + q - 1 - \sigma(b) = p + q - 1$$

$$2777 \quad 2778 \quad 2779 \quad 2780 \quad 2781 \quad 2782 \quad f_2 = \sigma(b) - p$$

2777 Hence, on computing the Expectation:

$$2778 \quad 2779 \quad 2780 \quad 2781 \quad 2782 \quad \mathbb{E}_X [f_1 X + f_2] = (p + q - 1) \mathbb{E}_X [X] + \sigma(b) - p = 0$$

$$2783 \quad 2784 \quad 2785 \quad 2786 \quad 2787 \quad 2788 \quad = (p + q - 1) \pi_1 + \sigma(b) - p = 0$$

2783 On simplifying:

$$2784 \quad 2785 \quad 2786 \quad 2787 \quad 2788 \quad 2789 \quad 2790 \quad \sigma(b) = \frac{p}{p+q}$$

$$2791 \quad 2792 \quad 2793 \quad 2794 \quad \implies b = \log \frac{p}{q}$$

2792 We could attempt to understand the characterization of the points on this manifold through local perturbation analysis.
2793 However, in this work, we classify them as stationary points and leave the comprehensive characterization for future research.
2794 Therefore,

2795 $\Gamma_{\text{station}}(p, q) \subseteq \mathbb{R}^4$ as follows:

$$2796 \quad 2797 \quad 2798 \quad 2799 \quad 2800 \quad 2801 \quad \Gamma_{\text{station}}(p, q) \triangleq \left\{ \gamma_{\min} = (e, w, b, a) \in \mathbb{R}^4 : e \neq 0, 1 + ae^2 = 0, 1 + 2w|w| = 0, b = \log \frac{p}{q} \right\},$$

2802 M.3.3. COMPUTING THE OPTIMAL BIAS

2803 Redefining quantities:

2804

$$\begin{aligned}
 f_1 &= \sigma \left(\underbrace{e^2 \left[\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2)| \right]}_{z_1} + b \right) \\
 &\quad - \sigma \left(\underbrace{e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2)| \right]}_{z_2} + b \right) \\
 &\quad + p + q - 1 \\
 f_2 &= \sigma \left(\underbrace{e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2)| \right]}_{z_2} + b \right) - p
 \end{aligned}$$

Finding the manifold of stationary points for the bias:

$$\begin{aligned}
 \frac{L(\theta)}{\partial b} &= \mathbb{E}_X [f_1 X + f_2] = 0 \\
 &= (\sigma(z_1) - \sigma(z_2) + p + q - 1) \mathbb{E}_X [X] + \sigma(z_2) - p = 0 \\
 &= (\sigma(z_1) - \sigma(z_2) + p + q - 1)\pi_1 + \sigma(z_2) - p = 0
 \end{aligned}$$

Hence on grouping terms together and simplifying:

$$\begin{aligned}
 &(\sigma(z_1) - \sigma(z_2) + p + q - 1)\pi_1 = p - \sigma(z_2) \\
 \implies &(\sigma(z_1) - 1)\pi_1 - \sigma(z_2)\pi_1 + p = p - \sigma(z_2) \\
 \implies &(\sigma(z_1) - 1)\pi_1 = \sigma(z_2)(\pi_1 - 1) \\
 \implies &\frac{\sigma(z_2)}{1 - \sigma(z_1)} = \frac{\pi_1}{1 - \pi_1} = \frac{p}{q}
 \end{aligned}$$

On using the definition of the sigmoid function and rearranging:

$$\begin{aligned}
 &\frac{1 + \exp(z_1)}{1 + \exp(-z_2)} = \frac{p}{q} \\
 \implies &\exp(z_1) + 1 = \frac{p}{q} (1 + \exp(-z_2)) \\
 \implies &\exp(2z_1) + \exp(z_1) = \frac{p}{q} \exp(z_1) + \frac{p}{q} \exp(z_1 - z_2) \\
 \implies &\exp(z_1)^2 + \exp(z_1) \left(1 - \frac{p}{q}\right) - \frac{p}{q} \cdot \exp(z_1 - z_2) = 0
 \end{aligned}$$

Basis the definitions of z_1 and z_2 , we have $z_1 - z_2 = e^2(1 + ae^2)(1 + 2w|w|)$. We define $A = \exp(e^2(1 + ae^2)(1 + 2w|w|))$. Hence we obtain a quadratic equations for for $\exp(z_1)$:

$$\exp(z_1)^2 + \exp(z_1) \left(1 - \frac{p}{q}\right) - \frac{p}{q} \cdot A = 0$$

On solving the quadratic equation for $\exp(z_1)$:

$$\exp(z_1) = \frac{1}{2} \left[\frac{p}{q} - 1 + \sqrt{\left(\frac{p}{q} - 1\right)^2 + 4 \cdot \frac{p}{q} \cdot A} \right] \quad (98)$$

$$\implies z_1 = \log \left(\frac{1}{2} \left[\frac{p}{q} - 1 + \sqrt{\left(\frac{p}{q} - 1\right)^2 + 4 \cdot \frac{p}{q} \cdot A} \right] \right) \quad (99)$$

$$\implies b_{\text{station}} = \log \left(\frac{1}{2} \left[\frac{p}{q} - 1 + \sqrt{\left(\frac{p}{q} - 1\right)^2 + 4 \cdot \frac{p}{q} \cdot A} \right] \right) \quad (100)$$

$$- e^2 \left[\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right], \quad (101)$$

where b_{station} denotes points where the bias remains a stationary quantity.

M.4. Proof of 15

Proof. Recall from Lemma 13 that for $\theta = (e, w, a) \in \mathbb{R}^3$, we have

$$\begin{aligned} \frac{\partial L}{\partial e} &= -\mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \cdot 2e (1 + ae^2) (1 + 2w|w|) \\ &\quad - \mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \cdot 2e^3 a (1 + 2w|w|), \\ \frac{\partial L}{\partial w} &= -\mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \cdot 2e^2 (1 + ae^2) (|w| + \text{sign}(w) w), \\ \frac{\partial L}{\partial a} &= -\mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \cdot e^4 (1 + 2w|w|), \end{aligned}$$

where $X \in \{0, 1\}$ is a Bernoulli random variable with $X \sim \text{Bern}(p/(p+q))$, and

$$\begin{aligned} f_1 &\triangleq 1 - p - q - \phi_1 + \phi_0, \quad f_2 \triangleq p - \phi_0, \\ \phi_1 &\triangleq \sigma \left(e^2 \left(\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right) + b_\star \right), \\ \phi_0 &\triangleq \sigma \left(e^2 \left(\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2) \right) + b_\star \right), \end{aligned}$$

where the optimal bias b_\star is obtained by solving $\pi_1 f_1 + f_2 = 0$. Using the definition of the gradient flow that $\dot{\theta} = -\nabla L(\theta)$ for $\theta = \theta_t$, we have

$$\begin{aligned} \dot{w} &= -\frac{L(\theta)}{\partial w} = \mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] \cdot 2e^2 (1 + ae^2) (|w| + \text{sign}(w) w) \\ \implies \frac{\dot{w}}{e^2 (|w| + \text{sign}(w) w)} &= \mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 2 (1 + ae^2). \end{aligned} \quad (102)$$

Similarly for a ,

$$\begin{aligned} \dot{a} &= -\frac{L(\theta)}{\partial a} = \mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] e^4 (1 + 2w|w|) \\ \implies \frac{\dot{a}}{e^4} &= \mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] (1 + 2w|w|) \end{aligned} \quad (103)$$

2915 Likewise, for e :

$$\begin{aligned}
 2916 \quad \dot{e} &= -\frac{L(\boldsymbol{\theta})}{\partial e} = \mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 2(1 + ae^2) e(1 + 2w|w|) \\
 2917 & \\
 2918 & \\
 2919 & \quad + \mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] (1 + 2w|w|) 2e^3 a \\
 2920 & \\
 2921 &
 \end{aligned} \tag{104}$$

2922 By substituting the expressions of 102 and 103 into 104:

$$\begin{aligned}
 2923 \quad \dot{e} &= \mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 2(1 + ae^2) e(1 + 2w|w|) \\
 2924 & \\
 2925 & \quad \underbrace{\hspace{10em}}_{102} \\
 2926 & \\
 2927 & \quad + \mathbb{E} \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] (1 + 2w|w|) 2e^3 a \\
 2928 & \\
 2929 & \quad \underbrace{\hspace{10em}}_{103} \\
 2930 & \\
 2931 &
 \end{aligned} \tag{105}$$

2932 We obtain:

$$\begin{aligned}
 2933 \quad \dot{e} &= \frac{\dot{w}}{2e^2 (|w| + \text{sign}(w) w)} e(1 + 2w|w|) + \frac{\dot{a}}{e^4} 2e^3 a \\
 2934 & \\
 2935 & \\
 2936 & \\
 2937 & \\
 2938 &
 \end{aligned} \tag{106}$$

2939 On rearranging and simplifying:

$$\begin{aligned}
 2940 \quad e\dot{e} &= \frac{\dot{w}}{(|w| + \text{sign}(w) w)} (1 + 2w|w|) + \dot{a}2a \\
 2941 & \\
 2942 & \\
 2943 & \implies e\dot{e} = \frac{\dot{w}}{2(|w|)} (1 + 2w|w|) + 2a\dot{a} \\
 2944 & \\
 2945 & \\
 2946 &
 \end{aligned} \tag{107}$$

2947 Integrating the above equation on both sides:

$$\begin{aligned}
 2948 \quad \int e\dot{e} &= \int \frac{\dot{w}}{4(|w|)} (1 + 2w|w|) + \int 2a\dot{a} \\
 2949 & \\
 2950 & \implies \frac{e^2(t)}{2} = \frac{\text{sign}(w(t)) \cdot \log |w(t)| + w(t)^2}{2} + a(t)^2 + \frac{c}{2} \\
 2951 & \\
 2952 & \implies e^2(t) = \text{sign}(w(t)) \cdot \log |w(t)| + w(t)^2 + 2a(t)^2 + c \\
 2953 & \\
 2954 & \\
 2955 & \\
 2956 &
 \end{aligned} \tag{108}$$

2957 Note that here $c \in \mathbb{R}$, is a quantity that depends on the initial conditions. Thus the energy $\mathcal{E}(\boldsymbol{\theta}_t) = \mathcal{E}(\boldsymbol{\theta}_0)$ for $w_0 \neq 0$.

2958 \square

2960 M.5. Role of Standard Initialization

2961 *Proof.* [Informal] Reproducing the derivatives for $\boldsymbol{\theta} = (e) \in \mathbb{R}^1$, assuming optimal b :

$$\begin{aligned}
 2962 \quad \frac{L(\boldsymbol{\theta})}{\partial e} &= \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 2e(1 + ae^2)(1 + 2w|w|) \\
 2963 & \\
 2964 & \\
 2965 & \quad + \mathbb{E}_X \left[(f_1 X + f_2) \left(X - \frac{1}{2} \right) \right] 2e^3 a(1 + 2w|w|) = 0 \\
 2966 & \\
 2967 & \\
 2968 & \\
 2969 &
 \end{aligned} \tag{109}$$

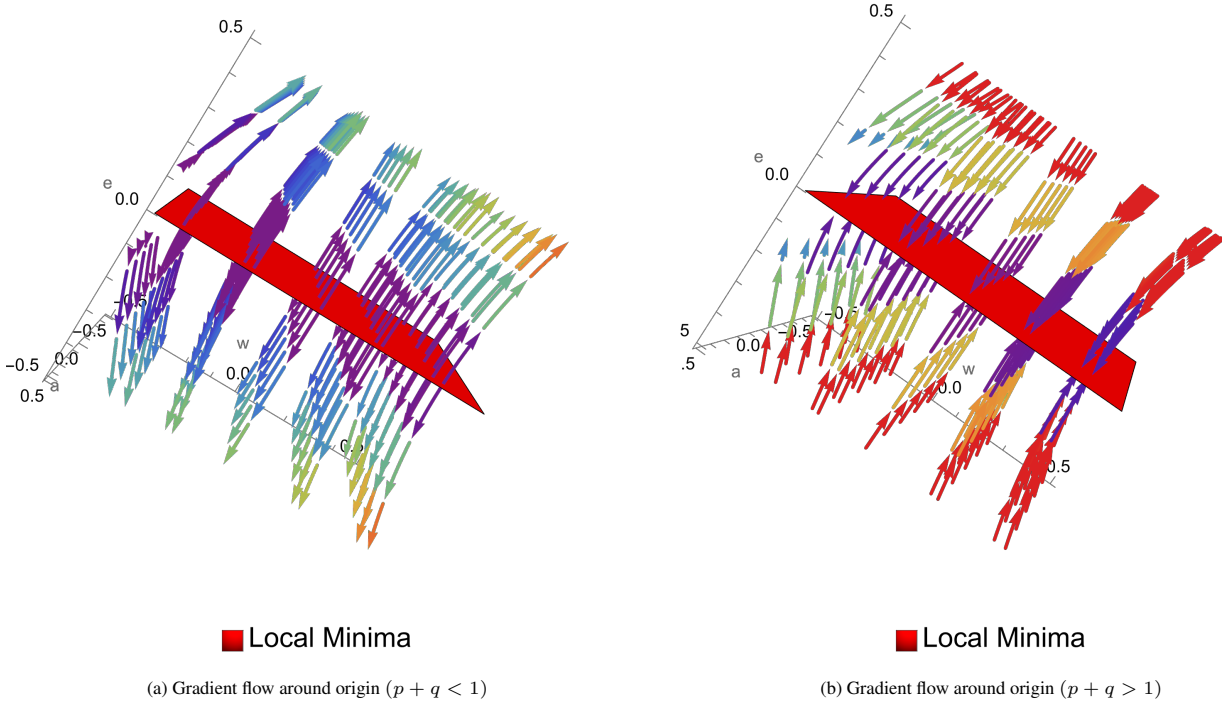


Figure 5: Gradient flow dynamics in \mathbb{R}^3 , near the origin, for the transformer parameters with attention scalar a (Sec. ??). The local minima are repellers for $p + q < 1$, while attracting for $p + q > 1$.

Where:

$$\begin{aligned}
 f_1 &= \sigma \left(\underbrace{e^2 \left[\frac{1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2)| \right]}_{z_1} + b \right) \\
 &\quad - \sigma \left(\underbrace{e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2)| \right]}_{z_2} + b \right) \\
 &\quad + p + q - 1 \\
 f_2 &= \sigma \left(\underbrace{e^2 \left[\frac{-1}{2} (1 + ae^2) (1 + 2w|w|) + w|w|(1 + ae^2)| \right]}_{z_2} + b \right) - p
 \end{aligned}$$

We assume the initialization is very small, making any product of quantities in $\theta = (e, w, a, b)$ much smaller than the individual quantities. Therefore, we can consider these products to be approximately zero. That is, $\forall x, y \in \theta, x \geq xy$ & $y \geq xy$ & $xy \approx 0$. Hence:

$$\begin{aligned}
 z_1 &= \sigma(b) \\
 z_2 &= \sigma(b)
 \end{aligned}$$

$$f_1 = p + q - 1$$

$$f_2 = \sigma(b) - p$$

Similarly:

$$\frac{L(\theta)}{\partial e} = 2\mathbb{E}_X [(f_1 X + f_2)(X)] e \tag{110}$$

Analysing $\mathbb{E}_X [(f_1 X + f_2)(X)]$:

$$\begin{aligned} \mathbb{E}_X [(f_1 X + f_2) X] &= (f_1 + f_2)\pi_1 \\ &= f_1\pi_1 - f_1\pi_1^2 \\ &= (p + q - 1)(\pi_1 - \pi_1^2) \end{aligned}$$

We used the fact that b is optimal in the above equations, specifically where $f_1\pi_1 + f_2 = 0$. On computing the gradient flow:

$$\dot{e} = -\frac{L(\theta)}{\partial e} = -(p + q - 1)(\pi_1 - \pi_1^2)e \implies e = e_0 \exp(-(p + q - 1)(\pi_1 - \pi_1^2)t)$$

Since $(p + q - 1)(\pi_1 - \pi_1^2) > 0$, $e \rightarrow 0$, which denotes it converges to the local minima.

□

N. Additional empirical results

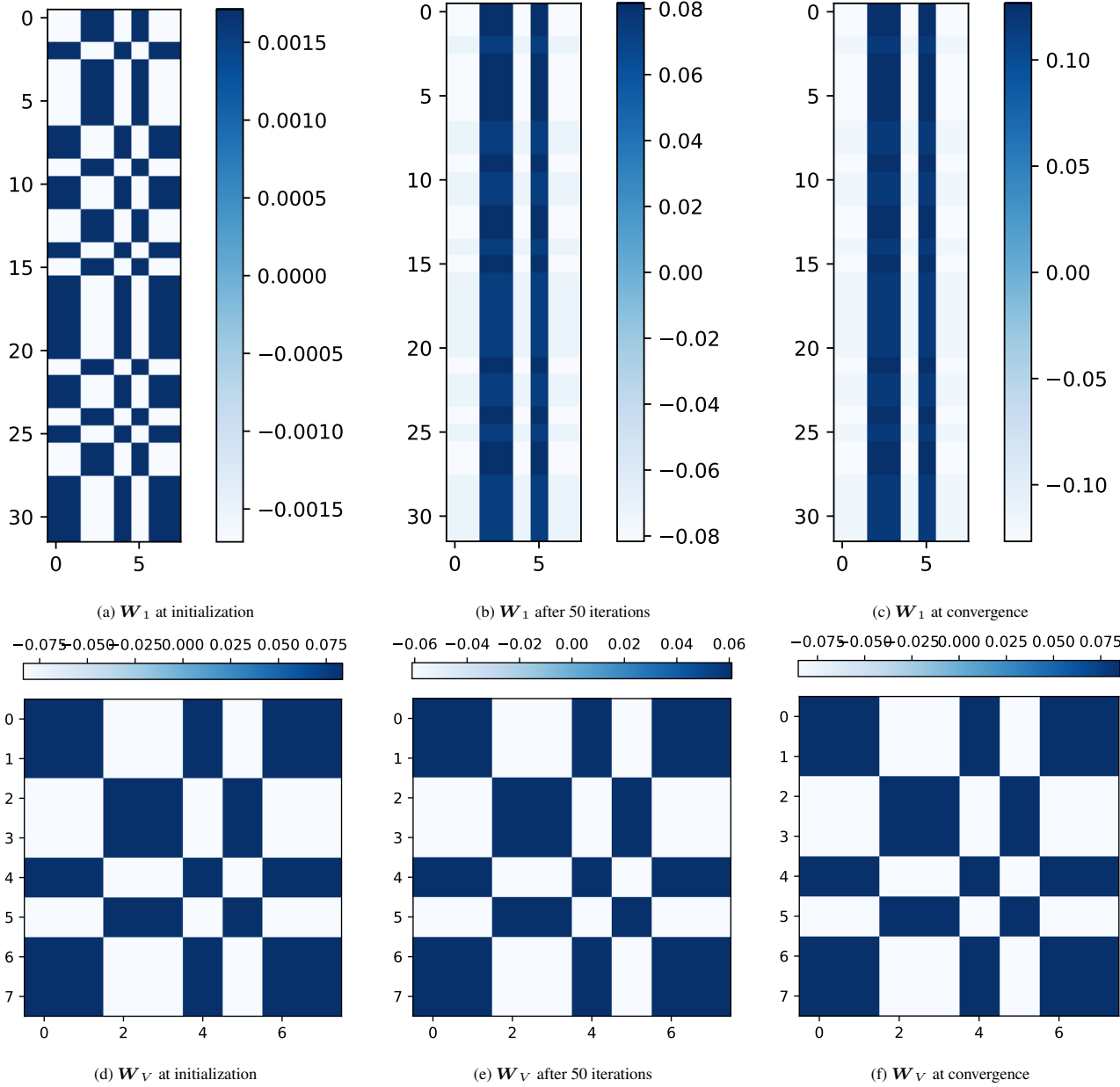


Figure 6: Evolution of parameters W_1 and W_V across iterations, starting from a rank-one initialization. The parameters maintain a rank-one structure across the entire training.

O. 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	$d \times 1$
transformer.h.attn.c_attn	$3d \times d$
transformer.h.attn.c_proj	$d \times d$
transformer.h.ln_2	$d \times 1$
transformer.h.mlp.c_fc	$4d \times d$
transformer.h.mlp.c_proj	$d \times 4d$
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	1
Attention heads	1
Repetitions	3 or 5